

SPATIO-TEMPORAL CHANGE-POINT DETECTION AND CONSTRAINED BAYESIAN OPTIMIZATION

A Thesis
Presented to
The Academic Faculty

by

Junzhuo Chen

In Partial Fulfillment
of the Requirements for the Degree
Doctor of Philosophy in the
H. Milton Stewart School of Industrial and Systems Engineering

Georgia Institute of Technology
May 2019

Copyright © 2019 by Junzhuo Chen

SPATIO-TEMPORAL CHANGE-POINT DETECTION AND CONSTRAINED BAYESIAN OPTIMIZATION

Approved by:

Professor Seong-Hee Kim, Advisor
H. Milton Stewart School of Industrial
and Systems Engineering
Georgia Institute of Technology

Professor Yao Xie, Advisor
H. Milton Stewart School of Industrial
and Systems Engineering
Georgia Institute of Technology

Professor Mustafa M. Aral
Department of Civil Engineering
Bartın University

Professor Kamran Paynabar
H. Milton Stewart School of Industrial
and Systems Engineering
Georgia Institute of Technology

Professor Jianjun Shi
H. Milton Stewart School of Industrial
and Systems Engineering
Georgia Institute of Technology

Date Approved: March 15, 2019

To my family, and all my friends.

ACKNOWLEDGEMENTS

First and foremost, I would like to express my deepest gratitude to my advisors, Dr. Seong-Hee Kim and Dr. Yao Xie, for their support, encouragement and trust during these years. Their guidance helped me in all the time of research and writing of this thesis. From them, I learned how to become a researcher.

I would like to thank my thesis committee members: Dr. Mustafa M. Aral, Dr. Jianjun Shi and Dr. Kamran Paynabar for their kindness in evaluating my thesis and their valuable feedback. I would also like to thank Dr. Craig Tovey, Dr. Robert Foley, Dr. David Goldberg, Dr. Nagi Gebraeel, Dr. Alexander Shapiro and Dr. Yajun Mei for their excellent courses. In addition, I would like to thank Amanda Ford and Dr. Dawn Strickland for their dedicated services at ISyE.

I would like to thank all my fellow students and friends accompanying me during my study at Tech. The five years that I spent with them in Atlanta is so wonderful and will be such a precious memory in my life.

Last but not least, I would like to thank my parents for their unconditional love and support.

TABLE OF CONTENTS

DEDICATION	iii
ACKNOWLEDGEMENTS	iv
LIST OF TABLES	viii
LIST OF FIGURES	ix
SUMMARY	xi
I INTRODUCTION	1
1.1 Literature Review	3
1.1.1 Spatio-temporal change-point detection	3
1.1.2 Optimal sensor network design	5
II TO REDUCE OR NOT TO REDUCE: A STUDY ON SPATIO-TEMPORAL CHANGE-POINT DETECTION	7
2.1 Background	9
2.1.1 Notation and Problem	9
2.1.2 LR-F-MCUSUM chart	10
2.1.3 T^2 -F-MCUSUM chart	11
2.2 RD-MCUSUM Charts	12
2.2.1 LR-RD-MCUSUM chart	12
2.2.2 T^2 -RD-MCUSUM chart	13
2.3 Theoretical Analysis for Effects of Spatial Correlation	13
2.3.1 Relations between statistics in full and reduced-dimension charts	14
2.3.2 Performance metric: ARL_1 measure	15
2.4 Experiments	20
2.4.1 Experimental setup	20
2.4.2 Results	20
2.5 Application: Water Quality Monitoring	23
2.5.1 Data	23
2.5.2 Spatial Models	24
2.5.3 Results	27

2.6	Conclusion	29
III	S^3T: A SCORE STATISTIC FOR SPATIO-TEMPORAL CHANGE-POINT DETECTION	30
3.1	Problem Formulation	31
3.2	Statistic for Offline and Online Detection	35
3.2.1	Quadratic score statistic	36
3.2.2	S^3T statistic for offline change-point detection	37
3.2.3	S^3T statistic for online change-point detection	37
3.3	Theoretical Approximations	38
3.3.1	Significance level for offline S^3T statistic	38
3.3.2	In-control Average Run Length (ARL_0) for online S^3T statistic . .	41
3.4	Numerical Examples	43
3.4.1	Simulation	43
3.4.2	Real data example: Solar flare detection	44
3.4.3	Case study: Water quality monitoring	45
3.5	Conclusions	47
IV	COMBINING CONSTRAINED BAYESIAN OPTIMIZATION AND SPATIO-TEMPORAL CHANGE-POINT DETECTION FOR SENSOR NETWORK DESIGN	48
4.1	Background	49
4.1.1	Constrained black-box function optimization	49
4.1.2	Bayesian Optimization	50
4.2	Confidence-Set based Constrained Bayesian Optimization	52
4.2.1	Connection to confidence bounds	54
4.2.2	Choices for h_1 and h_2	55
4.2.3	Extend to multi-task Gaussian processes (MTGP)	57
4.3	Combine CSCBO and S^3T for Sensor Network Design	58
4.3.1	Formulation	59
4.3.2	Measurement error	60
4.3.3	Process simulation	63
4.3.4	Wasserstein similarity metric	63
4.4	Experiments	65

4.4.1	Simulation setup	65
4.4.2	Error-free sensors	65
4.4.3	Sensors with measurement error	68
4.5	Conclusion	71
V	DISCUSSION AND CONCLUSIONS	73
APPENDIX A	— APPENDIX FOR CHAPTER 2	74
APPENDIX B	— APPENDIX FOR CHAPTER 3	79
REFERENCES	89

LIST OF TABLES

1	Summary of charts with the full and RD approaches.	13
2	Drift and variance parameters of the LR based charts.	16
3	Drift and variance parameters of the T^2 -statistic based charts.	17
4	Detection performance (ARL_1) obtained using the non-overlapped (NOV) and overlapped (OV) sets of scan clusters. Numbers in parentheses are standard errors.	28
5	Detection delay comparison between LR-F-MCUSUM and LR-RD-MCUSUM charts.	29
6	Notations.	35
7	Simulated and approximated significance level when the signal $\{\mathbf{x}_\ell\}$ follows a VAR(1) model.	40
8	Simulated expected detection delay.	44
9	Simulated expected detection delay in hours (numbers in parentheses are standard errors).	47
10	Performance metrics of the optimal feasible solutions found by CSCBO and NP + PFM. The unit of ECEDD is hour.	67
11	ECEDD (in hours) of sensor placements marked by circles in Figure 18(a) and (c) using different detection statistics. Numbers in parentheses are standard errors.	69
12	Simulated ARL_0 using S^3T when no contaminant event. $ARL_{\text{target}} = 10000$	70
13	Performance metrics of the optimal feasible solution found by CSCBO with S^3T as the detection statistic. The unit of ECEDD is hour.	71

LIST OF FIGURES

1	(a) A monitored area with $p = 7 \times 7$ locations or sensors and illustration of the spatial scanning using a circular shaped region; (b) Mapping of a full-dimensional observation vector into reduced-dimensional vectors corresponding to scanning regions.	8
2	Example with $p = 5$, $\tilde{p} = 2$, and $\mu_{c,r} = [1, 1, 0, 0, 0]'$: (a) m_{LR}/\tilde{m}_{LR} as a function of ρ ; and (b) m_{T^2}/\tilde{m}_{T^2} as a function of ρ	18
3	Simulated ARL_1 of LR-F-MCUSUM and LR-RD-MCUSUM charts with $r_{out} = \sqrt{2}$: (a) spherical model, (b) polynomial model and (c) Matérn model. . . .	21
4	Simulated ARL_1 of T^2 -F-MCUSUM and T^2 -RD-MCUSUM charts with $r_{out} = \sqrt{2}$: (a) spherical model, (b) polynomial model and (c) Matérn model. . . .	22
5	Shape of the Altamaha River ([65]).	24
6	A stream network example with nine stream segments ($i = 1, \dots, 9$) and three monitoring locations s_1, s_2, s_3	25
7	Visualization of the spatial covariance matrix for the Altamaha River. Each block in the covariance matrix corresponds to a branch of the river with a matching color.	26
8	Two sets of scan clusters for spatial scanning: (a) non-overlapped clusters; (b) overlapped clusters. Red stars represent possible spill locations.	27
9	Diagram showing the concatenation of samples.	34
10	Sliding window of length w for online detection.	38
11	Histograms and q-q plots of $W(\theta, \tau)$ for fixed values of θ and τ : $\tau = 30$, $\theta = 0.3$ for (a) and (c); $\tau = 40$, $\theta = 0.2$ for (b) and (d).	41
12	Comparison of approximated and simulated ARL for (a) $p = 1$, (b) $p = 2$, and (c) $p = 9$	42
13	Detection of solar flare at $t = 227$: (left) snapshot of the original SDO data at $t = 227$; (right) overlapping image patches for dimensionality reduction.	45
14	Detection statistics on logarithmic scale.	46
15	Mean utility gap between the true optimal feasible solution and the solution found by CSCBO: (a) deterministic, $H = -0.8$, (b) deterministic, $H = -0.99$, (c) stochastic, $H = -0.8$ and (d) stochastic, $H = -0.99$	56
16	Comparison on log median utility gap between CSCBO with independent GPs (idgp) and multi-task GPs (mtgp).	57
17	Three different sensor placements (number of sensors $M = 2$) on a hypothetical river with $D = 6$ nodes. Sensors are marked by the red crosses. The stream distances between each node are also marked on the plots.	65

18	Optimal feasible solutions found by CSCBO (circle) and NP + PFM (triangle): (a) $M = 5$ and $b = 0.05$, (b) $M = 5$ and $b = 0.0001$, (c) $M = 7$ and $b = 0.05$ and (d) $M = 7$ and $b = 0.0001$	68
19	Optimal feasible solutions found by CSCBO with S^3T as the detection statistic: (a) $M = 5$ and $\zeta_1 = 0.01$, (b) $M = 5$ and $\zeta_1 = 0.005$, (c) $M = 7$ and $\zeta_1 = 0.01$ and (d) $M = 7$ and $\zeta_1 = 0.005$	72
20	Ratio of simulated ARL_1 's (blue) and ratio of ARL_1 measures (red) in the known shift cluster case.	76
21	Performance loss based on ARL_1 measure and simulated ARL_1	77

SUMMARY

This thesis makes contributions to two research topics: spatio-temporal change-point detection and constrained Bayesian optimization. Spatio-temporal change-point detection is concerned with detecting statistical anomalies based on multiple data streams collected at different locations. In Chapter 2 and Chapter 3, we address two challenges in spatio-temporal change-point detection: (i) how to deal with data with high dimensionality, and (ii) how to capture spatial and temporal correlations. Bayesian optimization is a prevalent approach for optimization problems defined by expensive-to-evaluate black-box functions. In Chapter 4, we develop a practical algorithm for optimization problems with black-box objective function and constraints.

In Chapter 2, we study dimension reduction via spatial scanning. The majority of control charts using scan statistics for spatio-temporal change-point detection use full observation vectors. To deal with high dimensionality, most of the dimension reduction techniques are done as a post-processing step rather than in the data acquisition stage and thus the full sample covariance matrix is required. In a high dimensional application, (i) the sample covariance matrix tends to be ill-conditioned due to a limited number of samples; (ii) inversion of such a sample covariance matrix causes numerical issues; (iii) aggregating information from all variables may lead to high communication costs in sensor networks. We consider a set of reduced-dimension (RD) control charts which perform dimension reduction during data acquisition by spatial scanning and avoid the computational difficulties and possibly high communication costs. We characterize the performance difference between the RD and the full observation approaches, under several common spatial correlation models, in terms of average run lengths. Our results show that the RD approach has little performance loss under the correlation models considered in this chapter while enjoying all the implementation benefits. Our theoretical analysis is verified by extensive numerical studies including water quality monitoring.

In Chapter 3, we propose an efficient score statistic, called the S^3T statistic, to detect the emergence of a spatially and temporally correlated signal from either fixed-sample or sequential data. The signal may cause a mean shift and/or a change in the covariance structure. The score statistic can capture both the spatial and temporal structures of the change, and hence, is particularly powerful in detecting weak signals. The score statistic is computationally efficient and statistically powerful. Our main theoretical contribution is analytical approximations of the false alarm rates of the detection procedures. Numerical experiments on simulated and real data, as well as a real case study of water quality monitoring, demonstrate the good performance of our procedure.

In Chapter 4, we study the problem of optimal sensor network design, which is formulated as a joint problem of constrained black-box function optimization and spatio-temporal change-point detection. We propose a practical algorithm called the Confidence-Set based Constrained Bayesian Optimization (CSCBO), which provides a flexible framework to handle noisy black-box function constraints and is easy to implement. We also extend the algorithm to tackle with a challenge that arises specifically in the sensor network design problem: we use the Wasserstein similarity metric to deal with high-dimensional binary decision variables. Finally, the S^3T statistic proposed in Chapter 3 is combined with CSCBO to identify optimal sensor network designs that are robust to sensor measurement errors.

CHAPTER I

INTRODUCTION

The rapid development of sensor technology and communication network has enabled on-line monitoring of statistical anomalies based on high volume spatio-temporal data in a variety of industrial and service systems. The abrupt emergence of such an anomaly will change the distribution of the data and will usually cause destructive consequences, and hence quick detection is desired. Examples of spatio-temporal change-point detection include disease outbreak detection [30], water quality monitoring [1], and computer network intrusion detection [42]. While many methods have been established for classic problems, developing efficient detection procedure for spatio-temporal data involves new challenges. The first issue centers on the high-dimensionality of the data. For example, in a sensor network, the number of sensors deployed can be as large as thousands, which incurs high implementation cost if the data streams need to be processed centrally. Another issue is the complicated structure of the spatial and temporal correlations of the data. Usually the underlying process measured by the sensors is a spatio-temporal process where correlations exist in both space and time. Hence, the capacity of capturing correlation information is critical for efficient change-point detection procedures. In this thesis, we address these challenges through the following methods:

- We develop reduced-dimension methods via spatial scanning. Dimension reduction is achieved by breaking the monitored area into clusters and constructing local statistics for each cluster. The reduced-dimension methods with spatial scanning enjoy the following computational benefits: (i) it avoids estimating the entire sample covariance matrix which is difficult if we have a limited number of samples, (ii) it avoids the inversion of a large covariance matrix which is likely to cause numerical issues. On the other hand, the spatial scanning scheme also enables distributed implementation since data only needs to be processed locally instead of centrally. This will largely

reduce communication cost among sensors if the sensor network is deployed on a large area. Our theoretical and numerical studies show that the RD approach has little performance loss comparing to methods with full observation vectors and hence should be a preferable method in practice.

- We develop a new score statistic, called the S^3T statistic, which captures both spatial and temporal correlation of a signal, and hence, is particularly powerful in detecting weak signals. The S^3T statistic also avoids inversion of the spatial-temporal covariance matrix and is hence computational efficient comparing to the maximum likelihood ratio statistic. We also develop theoretical approximations of the false alarm rates of the proposed detection procedures. Performance of the new methods is demonstrated via simulated data and real data.

In this thesis, we also study the problem of optimal sensor network design, which consists of two components: selection of sensor locations and design of detection statistics. The sensor network we study is used for water quality monitoring, and hence, the design of such a network should answer not only where the sensors should be placed in space but also how the data collected should be processed and analyzed. The objective is to achieve minimal detection delay when a contamination event occurs, and meanwhile, the sensor network is subject to the constraints on the probability of detection and the false alarm rate. We formulate the problem as a joint problem of constrained black-box function optimization and spatio-temporal change-point detection. Black-box function optimization refers to the optimization problems defined by functions that do not have analytical forms and usually can only be evaluated via computer simulations. Bayesian optimization (BO) is a prevalent method for unconstrained black-box function optimization problems. In Chapter 4, we develop a practical algorithm called the Confidence-Set based Constrained Bayesian Optimization (CSCBO), which can be used to solve problems with black-box function constraints. We also use the Wasserstein similarity metric to tackle with high-dimensional binary decision variables, which is a challenge that arises specifically in the sensor network design problem. Finally, we combine the S^3T statistic with CSCBO to identify optimal

sensor network designs that are robust to sensor measurement errors.

1.1 Literature Review

1.1.1 Spatio-temporal change-point detection

Multivariate control charts are commonly used for spatio-temporal change-point detection. From a stream of observation vectors, a control chart computes a sequence of monitoring statistics and detects a change whenever the monitoring statistic goes beyond certain control limits. The control limits are pre-specified according to the requirement for the in-control average run length (known as the ARL_0), which captures the false-alarm-rate of a control chart. Another related performance metric is the out-of-control average run length (known as the ARL_1), which represents the expected number of samples needed to raise an alarm. Our goal is to detect the change quickly, namely, to achieve a short detection delay or ARL_1 for a given targeted ARL_0 .

Classical multivariate control charts mainly consider low-dimensional problems with a few number of data streams. Commonly used multivariate control charts include the T^2 chart by [21], the multivariate exponentially weighted moving average (MEWMA) chart by [34], and the multivariate cumulative sum (MCUSUM) chart. The CUSUM chart [38] is widely adopted due to its good property in detecting small shifts and its efficient recursive implementation that facilitates online monitoring. In the multivariate setting, one may construct MCUSUM charts by directly using vector observations; such methods can be largely classified as the log likelihood ratio (LR) based and the Hotelling’s T^2 statistic based methods. Among them, [19] uses a log LR statistic assuming the vector observations are i.i.d. multivariate normal. [11] computes a T^2 statistic for each vector observation and then forms a CUSUM chart based on a sequence of T^2 statistics. A key difference between the LR based and the T^2 based MCUSUM statistics is that the former assumes known shift direction vectors, while the latter does not make such an assumption.

Compared with the settings of classical multivariate control charts, spatio-temporal change-point detection tends to be much higher dimensional. One important instance of spatial-temporal change-point detection is disease outbreak detection, where a decision

maker collects health measurements such as disease counts or mortality rate from a number of adjacent regions to monitor a potential public health hazard (see [68] for a review). For example, [61] consider a problem of simultaneously monitoring 200,000 indicators of excess mortality in the UK health system. In this setting, various MCUSUM charts have been developed that describes the spatial disease pattern in the entire monitored area. See, for example, [47] and [64]. Since a disease outbreak usually happens in a group of neighboring regions, [60] incorporates spatial scanning and defines a spatial cluster as a group of local regions lying within a circle of certain radii, which is an extended formulation of [48]. However, spatial correlation is not considered in [60]. Adopting a similar definition of spatial clusters in [60], [22] propose an MCUSUM chart, which constructs an LR statistic for each spatial cluster under the normal assumption and scans through all possible clusters to detect a possible outbreak. [30] revise this method by using an analytical formula developed by [25] to approximate the control limits. Then, [29] extend the previous work to more general distributions. Besides public health monitoring, [71] propose a generalized-likelihood-ratio-test control chart for monitoring a product surface data with the assumption that the surface data follow the multivariate normal distribution with an identity covariance matrix to detect a mean shift of a certain pattern. Another instance of high-dimensional spatial temporal change-point detection occurs in astronomical imaging, where high-resolution video streams are monitored for solar flare detection, as studied in [72] and [33]. In this problem, each observation is a 67,744 dimensional vector consisting of image pixels and the goal is to detect an emergence of a sparse signal. Sequential detection of a sparse change is an active area of research, see, e.g., the work by [74] and [32]. However, these work usually assumes independent data streams without considering spatial correlation. See, for example, [12] for the importance of capturing spatial correlation for effective monitoring.

To tackle the difficulties caused by high-dimensionality, a viable solution is to perform a dimensionality reduction technique. The existing dimension reduction techniques include the principal component analysis (PCA) such as [35]; random linear projections such as [51], [3] and [58]; and wavelet transform methods such as [28] and [70]. In addition, [18] propose a projection based method combined with a T^2 control chart, and [6] and [77] propose

a Bayesian hierarchical approach to dimension-reduced spatio-temporal modeling for the task of predicting a high-dimensional response from a high-dimensional predictor. These methods still require full observation vectors because dimension-reduction is performed as a *post-processing* step rather than *in the data acquisition stage*. In other words, while monitoring a process, full observation vectors need to be collected and thus high communication cost is still incurred in addition to a possibly ill-conditioned full sample covariance matrix.

1.1.2 Optimal sensor network design

The optimal selection of sensor locations for water quality monitoring network has been studied by many researchers. [63] provides a comprehensive review of past approaches. Among the recent works, [66] and [67] formulate the problem as a optimization problem with two objectives: minimizing the detection delay and maximizing the probability of detection. A genetic algorithm (GA) is used to solve the optimization problem. [40] and [41] formulate the problem as a constrained optimization problem to minimize the detection delay with a constraint on the probability of detection and adopt a combined procedure of the nested partition (NP) [52] and the penalty function with memory [39]. We refer to the combined procedure as NP + PFM, which is proved to converge almost surely to the true optimal feasible solution [39].

Optimal sensor placement is essentially a black-box function optimization problem as the objective function (detection delay) and the constraint (probability of detection) do not have analytical forms and can only be evaluated via simulations. Bayesian optimization (BO) [36] is a prevalent method for black-box function optimization problems. A BO algorithm typically models a function by a Gaussian process (GP) [45] and guide the search based on an acquisition function. [59] provides an overview on how BO algorithms are applied in various tasks. Recently, some researchers extend BO algorithms to optimization problems with black-box function constraints. These works include [14], [2], [20], [27] and [31]. Typically, applications of BO algorithms have continuous decision variable with relatively low dimensionality. However, in the problem of sensor network design, the decision variable is the set of sensor locations and hence is a vector of categorical variable. Thus,

a reasonable similarity metric over sets of locations is critical to apply BO algorithms for sensor network design. [15] propose to use the earth mover's distance as the similarity metric and apply a unconstrained BO algorithm to find the optimal locations for weather sensors in the United Kingdom.

CHAPTER II

TO REDUCE OR NOT TO REDUCE: A STUDY ON SPATIO-TEMPORAL CHANGE-POINT DETECTION

In spatio-temporal change-point detection, one monitors an area using data streams measured at different locations and aims to detect any changes as soon as they occur. Typically, the data streams are observations of a particular quality index, which are sequentially collected by either physical sensors or professionals from different locations in the monitored area. In most cases, the quality index at each location can be modeled as either a discrete (e.g., network intrusion counts, or mortality) or a continuous (e.g., disease incidence rate, or contaminant concentration) random variable.

For high-dimensional applications, all control charts reviewed in Section 1.1.1 use full observation vectors when constructing monitoring statistics. Even for charts with the spatial scanning approach ([22]; [30] and [29]), a full-size covariance matrix is involved in computing the detection statistic. In the high-dimensional setting in the presence of spatial correlation, several difficulties exist in applying these control charts that require a full-size covariance matrix to real scenarios: (i) the sample covariance matrix tends to be ill-conditioned due to the number of samples being relatively small compared to the dimension of the covariance matrix in the streaming setting; (ii) the ill-conditioned sample covariance matrix causes numerical issues due to the matrix inversion involved in computing the statistics; and (iii) communication cost can be high for distributed sensor networks, since using full observation vectors means that all sensors need to exchange information with each other [17].

Another way to tackle the high-dimensionality problem is to perform *reduced-dimension* spatial scanning. One first breaks the entire monitoring area into overlapping local clusters of certain radii and only uses a subset of sensors or locations within the clusters, as in Figure 1. Then, one constructs a control-chart for each local cluster and detects a change whenever any of the local clusters fires an alarm. By doing so, each control chart only

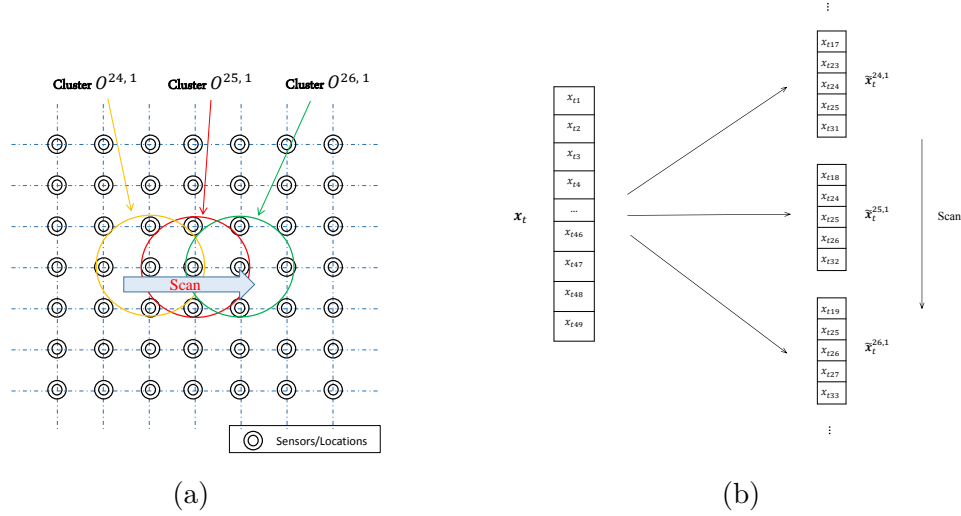


Figure 1: (a) A monitored area with $p = 7 \times 7$ locations or sensors and illustration of the spatial scanning using a circular shaped region; (b) Mapping of a full-dimensional observation vector into reduced-dimensional vectors corresponding to scanning regions.

monitors a small number of data streams that fall within the scanning cluster. In this chapter, we consider the reduced-dimensional spatial scanning with the spatial correlation *within the local cluster* only, rather than the full spatial correlation as in the earlier works such as [22]. Hence, our method, due to its spatial scanning nature, never needs to acquire full-dimensional observations and dimension reduction is performed *during data acquisition*. Our method is suitable for distributed processing required by sensor networks. A recent related work by [75] based on linear projection may also be used for spatial scanning; however, spatial correlation is not considered in that work. Moreover, due to the recursive nature of the CUSUM statistics, our method is suitable for *in situ* processing, which means that raw data are not needed to be stored, and this is again preferable in sensor networks.

Although the idea of the reduced-dimension spatial scanning is not new and, in fact, is often used in practice, one important question has never been unanswered in the literature: *How much do we lose by using reduced-dimension observations in the presence of spatial correlation?* We provide a precise answer to the amount of loss, by characterizing the difference of reduced-dimension charts and full-dimensional charts in term of their ARL_1 under a fixed ARL_0 . Our analysis shows that the RD approach usually spends $0 \sim 20\%$ more observations in ARL_1 than the full observation approach for a reasonable range of spatial

correlation among neighboring regions. Even when we lose more than 20%, the absolute differences in ARL_1 are small as one to five observations in most cases. In addition, we show that it is even possible that the RD approach may perform better than the full observation approach (i) when T^2 based charts for unknown shift directions are used or (ii) when some groups are completely independent of other groups as in a water quality monitoring example presented in Section 2.5. This is a blessing since an anomaly usually affects a local region and spatial correlation tends to decay with a distance. Thus, restricting to local sensors when forming a monitoring statistic should not significantly degrade the detection performance.

2.1 Background

In this section, we define notation and our problem. Then a few spatial scanning control charts are presented as representative charts that take the full observation approach.

2.1.1 Notation and Problem

In spatio-temporal change-point detection, observations (e.g., a quality index) are sequentially collected from different locations in the monitored area. Using the sequential observations, one desires to detect a possible change or anomaly in the monitored area as soon as possible.

Suppose there are p locations (sensors). For simplicity, we assume that the monitored area is rectangular and the locations (sensors) sit at a lattice of $p = MN$ points. This rectangular assumption on the shape of the monitored area can be relaxed, as we show in Section 2.5. Let $q_c = (m_c, n_c)$, where $m_c = 1, 2, \dots, M$ and $n_c = 1, 2, \dots, N$ denote the two-dimensional spatial coordinate of the location indexed by c , and $c = (m_c - 1)N + n_c$. Let $P = \{1, 2, \dots, p\}$ be the set of monitoring locations. At time t , the observation is a p -dimensional vector $\mathbf{x}_t = [x_{t1}, x_{t2}, \dots, x_{tp}]'$. Assume that different observation vectors are temporally independent but spatially correlated with covariance matrix $\mathbf{\Sigma}$. Further assume that the covariance matrix $\mathbf{\Sigma}$ of \mathbf{x}_t is known or can be estimated from data. The change only affects the mean and the covariance matrix remains the same. Under the hypothesis of no change, the observations $\mathbf{x}_1, \mathbf{x}_2, \dots$ are i.i.d. normally distributed with a mean vector $\boldsymbol{\mu}_0$ and a covariance matrix $\mathbf{\Sigma}$. Alternatively, there exists a change-point κ in

time, which represents an anomaly, and a subset of neighboring locations are *affected* by the change-point. For the locations affected by the change, the means of their observations are shifted, while observations from the unaffected locations keep the same distribution. This corresponds to a shift in the mean vector from $\boldsymbol{\mu}_0$ to some other vector $\boldsymbol{\mu}_1$. Without loss of generality, assume that the observation vectors have been standardized so that $\boldsymbol{\mu}_0 = \mathbf{0}$ and $[\boldsymbol{\Sigma}]_{i,i} = 1, \forall i = 1, \dots, p$, where $[\cdot]_{i,j}$ denotes the (i, j) th element of a matrix.

Due to spatial correlation, an anomaly often affects a cluster of neighboring locations. We assume that the shift cluster is circular shaped to facilitate the notation (but the cluster does not need to be circular shaped, which will be demonstrated in Section 2.5). A cluster is a set of locations $O^{c,r} = \{j \mid \|q_j - q_c\| \leq r, j \in P\}$, where c is the center of the cluster, r is the radius of the cluster, where $\|\cdot\|$ denotes the ℓ_2 -norm of a vector. In our setting (sensors are placed over a grid), the radius r is usually chosen from a discrete set of values. Let $R \subseteq \{1, \sqrt{2}, 2, 2\sqrt{2}, \dots\}$ be the set of possible values of r . Define an p -dimensional vector $[\boldsymbol{\mu}_{c,r}]_j = \delta_j$ for all $j \in O^{c,r}$ and 0, otherwise. Here, $[\cdot]_j$ denotes the j th element of a vector, and δ_j denotes the shift magnitude of the j th location. Hence, if an anomaly affects the cluster $O^{c,r}$, then $\boldsymbol{\mu}_1 = \boldsymbol{\mu}_{c,r}$.

Our goal is to detect anomalies that affect a cluster $O^{c,r}$, $c \in P$ and $r \in R$ by testing whether the mean of the observations has shifted from a nominal vector $\boldsymbol{\mu}_0$ to a different vector $\boldsymbol{\mu}_1$.

2.1.2 LR-F-MCUSUM chart

[22] consider a spatial scanning control chart based on LR statistics. Their method is based on full observation vectors and needs a full covariance matrix. Spatial scanning is achieved by zeroing out the part of the mean vector that falls out of the scanning region. Below, we review their method, referred to as the LR-F-MCUSUM chart hereafter. For a hypothetical shift cluster with center c and radius r , the sequence of monitoring statistics is given by,

$$S_t^{c,r} = \max_{1 \leq \tau \leq t} \sum_{i=\tau}^t \ell_i^{c,r}, \quad t = 1, 2, \dots, \quad (1)$$

where $\ell_i^{c,r}$ is the log LR statistic when the post-change mean vector is $\boldsymbol{\mu}_{c,r}$:

$$\ell_i^{c,r} = \boldsymbol{\mu}_{c,r}' \boldsymbol{\Sigma}^{-1} \left(\mathbf{x}_i - \frac{\boldsymbol{\mu}_{c,r}}{2} \right), \quad i = 1, 2, \dots, t. \quad (2)$$

When constructing the monitoring statistic in (1), we need to search over the unknown change-point time by maximizing with respect to τ , the variable that represents a putative change-point location. Note that the statistic for each cluster $S_t^{c,r}$ in (1) can be computed recursively

$$S_t^{c,r} = \max\{0, S_{t-1}^{c,r} + \ell_t^{c,r}\}, \quad t = 1, 2, \dots, \quad \text{and} \quad S_0^{c,r} = 0. \quad (3)$$

If the true sfhit center and radius (c, r) are known, we may perform the change-point detection when $S_t^{c,r}$ exceeds a pre-specified control limit. In practice, usually neither the shift center c nor the radius r is known a priori. In this case, one has to scan over all possible values of c and r , calculate the corresponding MCUSUM statistic, and form a global detection statistic by taking the maximum: $S_t^{**} = \max_{c \in P, r \in R} S_t^{c,r}$ for $t = 1, 2, \dots$. A change is detected whenever S_t^{**} exceeds a pre-specified control-limit $h_\ell^{**} > 0$, which is specified according to the requirement for the ARL_0 .

2.1.3 T^2 -F-MCUSUM chart

When there is no prior information on the magnitude or direction of the mean shift, a T^2 -statistic based MCUSUM chart is more appropriate. To have a reasonable comparison with reduced dimension charts, we introduce a T^2 -F-MCUSUM chart based on full dimensional observation vectors. It performs spatial scanning while using the full covariance matrix as analogous to [22]. Given a cluster $O^{c,r}$, we modify the full observation vector \mathbf{x}_t by replacing all elements that are not in the cluster with zeros $[\mathbf{x}_t^{c,r}]_j = [\mathbf{x}_t]_j$ for $j \in O^{c,r}$ and 0, otherwise. For each modified observation $\mathbf{x}_t^{c,r}$, compute a T^2 statistic

$$a_t^{c,r} = \mathbf{x}_t^{c,r}' \boldsymbol{\Sigma}^{-1} \mathbf{x}_t^{c,r} - \mu_T - k\sigma_T, \quad (4)$$

where k is a positive real-valued constant, $\mu_T = \text{E}[\mathbf{x}_t^{c,r}' \boldsymbol{\Sigma}^{-1} \mathbf{x}_t^{c,r}]$ and $\sigma_T^2 = \text{Var}[\mathbf{x}_t^{c,r}' \boldsymbol{\Sigma}^{-1} \mathbf{x}_t^{c,r}]$, which are the in-control mean and standard deviation of the T^2 statistic, respectively. The calculations for μ_T and σ_T^2 are discussed in Section 4. The monitoring statistic $T_t^{c,r}$ is

computed recursively over time

$$T_t^{c,r} = \max\{0, T_{t-1}^{c,r} + d_t^{c,r}\}, \quad t = 1, 2, \dots, \quad \text{and} \quad T_0^{c,r} = 0. \quad (5)$$

With unknown shift center and radius, we again search over all possible clusters and sizes to form the global detection statistic $T_t^{**} = \max_{c \in P, r \in R} T_t^{c,r}$ for $t = 1, 2, \dots$ and detection is performed by comparing T_t^{**} with a control limit h_a^{**} .

2.2 RD-MCUSUM Charts

In this section, we present the reduced-dimension approach. For each scan cluster, a control chart is constructed for reduced-dimension observation vectors while only considering local covariance. We develop two versions of MCUSUM charts, based on the LR statistic and the T^2 statistic, respectively.

2.2.1 LR-RD-MCUSUM chart

We start by considering an LR based chart, which is referred as the LR-RD-MCUSUM chart hereafter. For each scan cluster $O^{c,r}$, we truncate the original data vector \mathbf{x}_t into a lower dimensional vector $\tilde{\mathbf{x}}_t^{c,r}$, where $[\mathbf{x}_t]_j$ is positioned in $\tilde{\mathbf{x}}_t^{c,r}$ if $j \in O^{c,r}$, and is eliminated otherwise, as illustrated in Figure 1(b). The monitoring statistic for that particular cluster is computed over vectors of dimension $|O^{c,r}|$, where $|\cdot|$ denotes the cardinality of a set. At each time, the LR statistic for $O^{c,r}$ is computed

$$\tilde{\ell}_i^{c,r} = \tilde{\boldsymbol{\mu}}_{c,r}' \boldsymbol{\Sigma}_{c,r}^{-1} (\tilde{\mathbf{x}}_i^{c,r} - \frac{\tilde{\boldsymbol{\mu}}_{c,r}}{2}), \quad i = 1, 2, \dots, t. \quad (6)$$

Here $\tilde{\boldsymbol{\mu}}_{c,r}$ and $\boldsymbol{\Sigma}_{c,r}$ are the sub-vector and sub-matrix of $\boldsymbol{\mu}_{c,r}$ and $\boldsymbol{\Sigma}$, respectively. Then, based on (6), the detection statistic $\tilde{S}_t^{c,r}$ is computed recursively for each cluster similar to (3). Finally, a global monitoring statistic is formed by taking the maximum over all clusters and sizes

$$\tilde{S}_t^{**} = \max_{c \in P, r \in R} \tilde{S}_t^{c,r}, \quad t = 1, 2, \dots \quad (7)$$

An alarm is signaled whenever \tilde{S}_t^{**} exceeds a pre-specified control limit \tilde{h}_ℓ^{**} .

Table 1: Summary of charts with the full and RD approaches.

	F-MCUSUM	RD-MCUSUM
LR Based	$S_t^{**} = \max_{c,r} S_t^{c,r} \geq h_\ell^{**}$	$\tilde{S}_t^{**} = \max_{c,r} \tilde{S}_t^{c,r} \geq \tilde{h}_\ell^{**}$
T^2 Based	$T_t^{**} = \max_{c,r} T_t^{c,r} \geq h_a^{**}$	$\tilde{T}_t^{**} = \max_{c,r} \tilde{T}_t^{c,r} \geq \tilde{h}_a^{**}$

2.2.2 T^2 -RD-MCUSUM chart

Finally, similar to the reduced-dimension LR chart above, we construct the reduced-dimension T^2 statistic chart, which is referred to as the T^2 -RD-MCUSUM chart hereafter. Given a cluster, at each time, we compute a T^2 statistic,

$$\tilde{a}_t^{c,r} = \tilde{\mathbf{x}}_t^{c,r'} \Sigma_{c,r}^{-1} \tilde{\mathbf{x}}_t^{c,r} - \tilde{p} - k\sqrt{2\tilde{p}}, \quad i = 1, 2, \dots, t, \quad (8)$$

where $\tilde{p} = |O^{c,r}|$. Note that when the process is in-control, the mean and variance of T^2 statistic is given by

$$\mathbb{E}[\tilde{\mathbf{x}}_t^{c,r'} \Sigma_{c,r}^{-1} \tilde{\mathbf{x}}_t^{c,r}] = \tilde{p}, \quad \text{and} \quad \text{Var}[\tilde{\mathbf{x}}_t^{c,r'} \Sigma_{c,r}^{-1} \tilde{\mathbf{x}}_t^{c,r}] = 2\tilde{p}.$$

Based on $\tilde{a}_t^{c,r}$, the monitoring statistic for each cluster is computed recursively over time similar to (5). The global detection statistic is formed by maximizing over all clusters: $\tilde{T}_t^{**} = \max_{c \in P, r \in R} \tilde{T}_t^{c,r}$ for $t = 1, 2, \dots$. An alarm is signaled when \tilde{T}_t^{**} exceeds a control limit \tilde{h}_a^{**} .

Hereafter, the charts based on full observation vectors are referred to as the F-MCUSUM charts, including LR-F-MCUSUM and T^2 -F-MCUSUM charts, and the reduced dimension charts are referred to as the RD-MCUSUM charts, including LR-RD-MCUSUM and T^2 -RD-MCUSUM. Table 1 summarizes our methods and terminology.

2.3 Theoretical Analysis for Effects of Spatial Correlation

In this section, we compare analytically the performance of the F-MCUSUM charts with that of the RD-MCUSUM charts. We use ARL_1 given a target ARL_0 as our performance metric. To make a fair comparison among various control charts, we calibrate the control limits so that their actual ARL_0 are equal to the target values.

Both the F-MCUSUM and RD-MCUSUM charts use scan statistics to search for the true shift cluster from the set of all possible shift clusters $\{O^{c,r} | c \in P, r \in R\}$ at each time. In practice, the center and the radius of an anomaly are unknown so the entire monitored area needs to be scanned. In this section, we perform a theoretical analysis of the F-MCUSUM and RD-MCUSUM charts in a simplified setting, i.e., when the actual shift cluster is known. This simplified situation provides some insights into understanding the impact of dimensionality reduction on the performance of an MCUSUM chart with scanning.

Suppose a shift affects a cluster with center c and radius r , respectively, i.e., $O^{c,r}$ is the actual shift cluster. Suppose the cluster $O^{c,r}$ contains \tilde{p} locations. Without loss of generality, assume the affected locations correspond to the first \tilde{p} entries in the observation vectors, for instance, through reindexing. Hence, the post-change mean vector is $\boldsymbol{\mu}_1 = [\mu_1, \dots, \mu_{\tilde{p}}, 0, \dots, 0]' = [\tilde{\boldsymbol{\mu}}'_{c,r}, 0, \dots, 0]'$, where $\tilde{\boldsymbol{\mu}}'_{c,r} = [\mu_1, \dots, \mu_{\tilde{p}}]' \neq \mathbf{0}$. We have $\mathbf{x}_t = [\tilde{\mathbf{x}}_t^{c,r'}, 0, \dots, 0]' + [0, \dots, 0, \hat{\mathbf{x}}_t^{c,r'}]' = \mathbf{x}_t^{c,r} + [0, \dots, 0, \hat{\mathbf{x}}_t^{c,r'}]'$, where $\tilde{\mathbf{x}}_t^{c,r} = [x_{t1}, \dots, x_{t\tilde{p}}]'$ and $\hat{\mathbf{x}}_t^{c,r'} = [x_{t(\tilde{p}+1)}, \dots, x_{tp}]'$. Furthermore, partition the $p \times p$ dimensional covariance matrix for the full observation vectors accordingly $\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}$, where $\boldsymbol{\Sigma}_{11} \in R_{\tilde{p} \times \tilde{p}}$, $\boldsymbol{\Sigma}_{12} \in R_{\tilde{p} \times (p-\tilde{p})}$, $\boldsymbol{\Sigma}_{21} \in R_{(p-\tilde{p}) \times \tilde{p}}$ and $\boldsymbol{\Sigma}_{22} \in R_{(p-\tilde{p}) \times (p-\tilde{p})}$. Using the Schur complement [78] of $\boldsymbol{\Sigma}_{11}$, we write the inverse of $\boldsymbol{\Sigma}$ as

$$\boldsymbol{\Sigma}^{-1} = \begin{bmatrix} \boldsymbol{\Sigma}_{11}^{-1} + \boldsymbol{\Sigma}_{A^*}^{-1} & \boldsymbol{\Sigma}_B^{-1} \\ \boldsymbol{\Sigma}_C^{-1} & \boldsymbol{\Sigma}_D^{-1} \end{bmatrix}, \quad \text{where}$$

$$\begin{aligned} \boldsymbol{\Sigma}_{A^*}^{-1} &= \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12} (\boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12})^{-1} \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1}; & \boldsymbol{\Sigma}_B^{-1} &= -\boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12} (\boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12})^{-1}; \\ \boldsymbol{\Sigma}_C^{-1} &= -(\boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12})^{-1} \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1}; & \text{and} & \quad \boldsymbol{\Sigma}_D^{-1} = (\boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12})^{-1}. \end{aligned} \tag{9}$$

2.3.1 Relations between statistics in full and reduced-dimension charts

As a basis for the subsequent analysis, we derive relations for the LR and T^2 statistics used in the F- and the RD-MCUSUM charts. Recall that the likelihood ratio and the T^2 statistic based on full observations are $\ell_t^{c,r}$ and $a_t^{c,r}$, respectively and that their counterparts based

on reduced-dimension vectors are $\tilde{\ell}_t^{c,r}$ and $\tilde{a}_t^{c,r}$, respectively. Using equations (2), (6) and (9), the original LR statistic (2) can be written as

$$\begin{aligned}\ell_t^{c,r} &= \tilde{\mu}_{c,r}' \Sigma_{11}^{-1} \tilde{\mathbf{x}}_t^{c,r} - 0.5 \tilde{\mu}_{c,r}' \Sigma_{11}^{-1} \tilde{\mu}_{c,r} + \tilde{\mu}_{c,r}' \Sigma_{A^*}^{-1} \tilde{\mathbf{x}}_t^{c,r} + \tilde{\mu}_{c,r}' \Sigma_B^{-1} \hat{\mathbf{x}}_t^{c,r} - 0.5 \tilde{\mu}_{c,r}' \Sigma_{A^*}^{-1} \tilde{\mu}_{c,r} \\ &= \tilde{\ell}_t^{c,r} + \tilde{\mu}_{c,r}' \Sigma_{A^*}^{-1} \tilde{\mathbf{x}}_t^{c,r} + \tilde{\mu}_{c,r}' \Sigma_B^{-1} \hat{\mathbf{x}}_t^{c,r} - 0.5 \tilde{\mu}_{c,r}' \Sigma_{A^*}^{-1} \tilde{\mu}_{c,r}.\end{aligned}\quad (10)$$

From (10), we see that dimension reduction for LR based chart is equivalent to removing the last three terms on the right hand side of (10) from the full LR statistic. It is worth mentioning that $\tilde{\mu}_{c,r}' \Sigma_B^{-1} \hat{\mathbf{x}}_t^{c,r}$ contains only noise information since no mean shift occurs in $\hat{\mathbf{x}}_t^{c,r}$.

Similarly, we derive a relation for the T^2 statistics in the F- and the RD-MCUSUM charts. By equation (9) we obtain the in-control mean and variance of the full version T^2 statistic as follows:

$$\begin{aligned}\mathbb{E}[\mathbf{x}_t^{c,r'} \Sigma^{-1} \mathbf{x}_t^{c,r}] &= \tilde{p} + \text{tr}\{\Sigma_{A^*}^{-1} \Sigma_{11}\} \quad \text{and} \\ \text{Var}[\mathbf{x}_t^{c,r'} \Sigma^{-1} \mathbf{x}_t^{c,r}] &= 2\text{tr}\{\Sigma_{A^*}^{-1} \Sigma_{11} \Sigma_{A^*}^{-1} \Sigma_{11}\} + 4\text{tr}\{\Sigma_{A^*}^{-1} \Sigma_{11}\} + 2\tilde{p}.\end{aligned}$$

Hence, the statistics $a_t^{c,r}$ defined in (4) and $\tilde{a}_t^{c,r}$ defined in (8) are related to each other via

$$\begin{aligned}a_t^{c,r} &= \tilde{\mathbf{x}}_t^{c,r'} (\Sigma_{11}^{-1} + \Sigma_{A^*}^{-1}) \tilde{\mathbf{x}}_t^{c,r} - (\tilde{p} + \text{tr}\{\Sigma_{A^*}^{-1} \Sigma_{11}\}) - k[2\text{tr}\{\Sigma_{A^*}^{-1} \Sigma_{11} \Sigma_{A^*}^{-1} \Sigma_{11}\} + 4\text{tr}\{\Sigma_{A^*}^{-1} \Sigma_{11}\} + 2\tilde{p}]^{1/2} \\ &= \tilde{a}_t^{c,r} + \tilde{\mathbf{x}}_t^{c,r'} \Sigma_{A^*}^{-1} \tilde{\mathbf{x}}_t^{c,r} - \text{tr}\{\Sigma_{A^*}^{-1} \Sigma_{11}\} - k([2\text{tr}\{\Sigma_{A^*}^{-1} \Sigma_{11} \Sigma_{A^*}^{-1} \Sigma_{11}\} + 4\text{tr}\{\Sigma_{A^*}^{-1} \Sigma_{11}\} + 2\tilde{p})^{1/2} - (2\tilde{p})^{1/2}.\end{aligned}$$

2.3.2 Performance metric: ARL_1 measure

We aim to compare the detection performance of the F- and the RD-MCUSUM charts in terms of their ARL_1 for a fixed ARL_0 . To do so, we will define a performance metric called the ARL_1 measure. For a fixed ARL_0 , a smaller ARL_1 measure implies a smaller ARL_1 .

[25] derive a formula to approximate ARL for both in-control and out-of-control processes of a single CUSUM chart, which can be used to approximate ARL_1 for a fixed target value, ARL_0 . The formula is given by

$$\text{ARL} \approx \begin{cases} \frac{\Omega^2}{2d^2} \left\{ \exp \left[-\frac{2d(H + 1.166\Omega)}{\Omega^2} \right] - 1 + \frac{2d(H + 1.166\Omega)}{\Omega^2} \right\}, & \text{if } d \neq 0; \\ \left(\frac{H + 1.166\Omega}{\Omega} \right)^2, & \text{if } d = 0, \end{cases} \quad (11)$$

Table 2: Drift and variance parameters of the LR based charts.

	LR-F-MCUSUM	LR-RD-MCUSUM
d_0	$-\frac{1}{2}\boldsymbol{\mu}_{c,r}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_{c,r}$	$-\frac{1}{2}\tilde{\boldsymbol{\mu}}_{c,r}'\boldsymbol{\Sigma}_{11}^{-1}\tilde{\boldsymbol{\mu}}_{c,r}$
$d_{c,r}$	$\frac{1}{2}\boldsymbol{\mu}_{c,r}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_{c,r}$	$\frac{1}{2}\tilde{\boldsymbol{\mu}}_{c,r}'\boldsymbol{\Sigma}_{11}^{-1}\tilde{\boldsymbol{\mu}}_{c,r}$
$\Omega_0^2(\Omega_{c,r}^2)$	$\boldsymbol{\mu}_{c,r}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_{c,r}$	$\tilde{\boldsymbol{\mu}}_{c,r}'\boldsymbol{\Sigma}_{11}^{-1}\tilde{\boldsymbol{\mu}}_{c,r}$

where d is the drift parameter, Ω^2 is the variance parameter, and H is the control limit.

In our settings, specifically, if observations are temporally independent, for the likelihood ratio statistics based chart, $d = E[\ell_t^{c,r}]$, $\Omega^2 = \text{Var}[\ell_t^{c,r}]$ for the LR-F-MCUSUM chart, and $d = E[\tilde{\ell}_t^{c,r}]$, $\Omega^2 = \text{Var}[\tilde{\ell}_t^{c,r}]$ for the LR-RD-MCUSUM chart; in the T^2 based charts, $d = E[a_t^{c,r}]$, $\Omega^2 = \text{Var}[a_t^{c,r}]$ for the T^2 -F-MCUSUM chart and $d = E[\tilde{a}_t^{c,r}]$, $\Omega^2 = \text{Var}[\tilde{a}_t^{c,r}]$ for the T^2 -RD-MCUSUM chart. In the following, we denote the in-control drift and variance parameters as d_0 and Ω_0^2 , respectively. The out-of-control parameters are defined similarly. If the shift center is c and radius is r , we denote the out-of-control drift and variance as $d_{c,r}$ and $\Omega_{c,r}^2$, respectively. Tables 2 and 3 summarize these parameters.

The in-control drift d_0 is always negative but the out-of-control drift $d_{c,r}$ can be either negative (if the shift magnitude is too small) or positive. In this section, we assume that the shift magnitude is large enough so that $d_{c,r}$ is positive.

Consider a special function called the Lambert W function. Let $W_0(\cdot)$ be the principal branch of the Lambert W function [10]. For fixed target ARL_0 , $d_0 < 0$ and $d_{c,r} > 0$, using (11) and the notations defined above, we derive an approximation for ARL_1 as follows:

$$\text{ARL}_1 \approx -\frac{\Omega_0^2}{2d_0d_{c,r}}\epsilon_{\eta_0} - 1.166\frac{1}{d_{c,r}}(\Omega_0 - \Omega_{c,r}), \quad (12)$$

where

$$\epsilon_{\eta_0} = -W_{-1}(-e^{-\eta_0}) - \eta_0, \quad \text{and} \quad \eta_0 = \frac{2d_0^2}{\Omega_0^2}\text{ARL}_0 + 1.$$

The derivation is presented in Appendix A.1. It is noteworthy that the ϵ_{η_0} function tends to be very flat. Moreover, with a fixed ARL_0 , the values of η_0 for the F- and RD-MCUSUM charts are very close. In addition, for the LR based charts, since $\Omega_0 = \Omega_{c,r}$, the second term in the right hand side of (12) is equal to zero, and for the T^2 statistic based charts, it is small compared to the first term. Thus, when we fix ARL_0 to compare ARL_1 between the

Table 3: Drift and variance parameters of the T^2 -statistic based charts.

T^2 -F-MCUSUM	
d_0	$-k[2\text{tr}\{\Sigma_{A^*}^{-1}\Sigma_{11}\Sigma_{A^*}^{-1}\Sigma_{11}\} + 4\text{tr}\{\Sigma_{A^*}^{-1}\Sigma_{11}\} + 2\tilde{p}]^{1/2}$
$d_{c,r}$	$\tilde{\mu}'_{c,r}(\Sigma_{11}^{-1} + \Sigma_{A^*}^{-1})\tilde{\mu}_{c,r} - k[2\text{tr}\{\Sigma_{A^*}^{-1}\Sigma_{11}\Sigma_{A^*}^{-1}\Sigma_{11}\} + 4\text{tr}\{\Sigma_{A^*}^{-1}\Sigma_{11}\} + 2\tilde{p}]^{1/2}$
Ω_0^2	$2\text{tr}\{\Sigma_{A^*}^{-1}\Sigma_{11}\Sigma_{A^*}^{-1}\Sigma_{11}\} + 4\text{tr}\{\Sigma_{A^*}^{-1}\Sigma_{11}\} + 2\tilde{p}$
$\Omega_{c,r}^2$	$\Omega_0^2 + 4\tilde{\mu}'_{c,r}(\Sigma_{11}^{-1} + \Sigma_{A^*}^{-1})\Sigma_{11}(\Sigma_{11}^{-1} + \Sigma_{A^*}^{-1})\tilde{\mu}_{c,r}$
T^2 -RD-MCUSUM	
d_0	$-k(2\tilde{p})^{1/2}$
$d_{c,r}$	$\tilde{\mu}'_{c,r}\Sigma_{11}^{-1}\tilde{\mu}_{c,r} - k(2\tilde{p})^{1/2}$
Ω_0^2	$2\tilde{p}$
$\Omega_{c,r}^2$	$2\tilde{p} + 4\tilde{\mu}'_{c,r}\Sigma_{11}^{-1}\tilde{\mu}_{c,r}$

F- and RD-MCUSUM charts, we may compare the values of $|\Omega_0^2/(d_0 d_{c,r})|$, which we call the ARL_1 measure.

In the following, denote the ARL_1 measure for LR-F-MCUSUM and LR-RD-MCUSUM charts as m_{LR} and \tilde{m}_{LR} , respectively; and for T^2 -F-MCUSUM and T^2 -RD-MCUSUM charts as m_{T^2} and \tilde{m}_{T^2} , respectively. For LR based charts, we can obtain,

$$m_{LR} = \frac{4}{\tilde{\mu}'_{c,r}\Sigma^{-1}\tilde{\mu}_{c,r}} = \frac{4}{\tilde{\mu}'_{c,r}\Sigma_{11}^{-1}\tilde{\mu}_{c,r} + \tilde{\mu}'_{c,r}\Sigma_{A^*}^{-1}\tilde{\mu}_{c,r}} \leq \frac{4}{\tilde{\mu}'_{c,r}\Sigma_{11}^{-1}\tilde{\mu}_{c,r}} = \tilde{m}_{LR}. \quad (13)$$

Equation (13) shows that m_{LR} is always smaller than or equal to \tilde{m}_{LR} , which, in turn, implies that ARL_1 of the method using full observation vectors is always smaller than that of the chart with reduced dimension vectors. As a smaller ARL_1 measure implies a smaller ARL_1 , we expect that the LR-F-MCUSUM chart generally detects a shift faster than the LR-RD-MCUSUM chart.

For T^2 based charts, we have

$$m_{T^2} = \left[k^2 - k \frac{\tilde{\mu}'_{c,r}(\Sigma_{11}^{-1} + \Sigma_{A^*}^{-1})\tilde{\mu}_{c,r}}{[2\text{tr}\{\Sigma_{A^*}^{-1}\Sigma_{11}\Sigma_{A^*}^{-1}\Sigma_{11}\} + 4\text{tr}\{\Sigma_{A^*}^{-1}\Sigma_{11}\} + 2\tilde{p}]^{1/2}} \right]^{-1}$$

and

$$\tilde{m}_{T^2} = \left[k^2 - k \frac{\tilde{\mu}'_{c,r}\Sigma_{11}^{-1}\tilde{\mu}_{c,r}}{\sqrt{2\tilde{p}}} \right]^{-1}.$$

Note that the theoretical performance measure and the above analysis are applicable to spatial covariance matrix, Σ , with a general structure. Several commonly used spatial

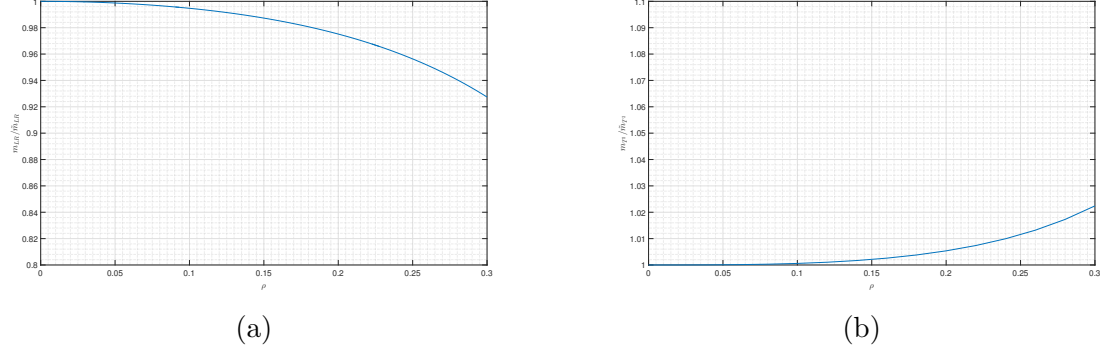


Figure 2: Example with $p = 5$, $\tilde{p} = 2$, and $\mu_{c,r} = [1, 1, 0, 0, 0]'$: (a) $m_{\text{LR}}/\tilde{m}_{\text{LR}}$ as a function of ρ ; and (b) m_{T^2}/\tilde{m}_{T^2} as a function of ρ .

covariance structures are as follows. In the following, d denotes the distance between two sensors, $C(d|\rho)$ denotes the correlation function between two sensors, which is a function of d and some parameters. Below, $\mathbb{1}_{\{\cdot\}}$ denotes the indicator function of an event, which takes value 1 only when the event is true.

1. Spherical model: $C(d|\rho) = \mathbb{1}_{\{d=0\}} + \rho\mathbb{1}_{\{d=1\}} + \frac{\rho}{2}\mathbb{1}_{\{d=\sqrt{2}\}}$ for $\rho \in [0, 1]$.
2. Polynomial model: $C(d|\rho) = \mathbb{1}_{\{d=0\}} + \rho^d\mathbb{1}_{\{d>0\}}$ for $\rho \in [0, 1]$.
3. Matérn model: $C(d|\theta) = \mathbb{1}_{\{d=0\}} + \frac{1}{2^{v-1}\Gamma(v)}(\sqrt{2}v^{1/2}d/\theta)^v K_v(\sqrt{2}v^{1/2}d/\theta)\mathbb{1}_{\{d>0\}}$ for $\theta > 0$ where K_v is the modified Bessel function of order v (See [46]).

Based on the correlation function, the entries of the covariance matrix $[\Sigma]_{i,j}$ is determined as $C(d(q_i, q_j)|\theta)$, where q_i and q_j are the coordinates of sensors i and j , respectively.

2.3.2.1 An illustrative example

Using a simple illustrative example, we calculate how much we lose or gain in terms of ARL_1 when reduced dimension vectors are used. Although the ARL_1 measure is applicable to general spatial correlation structure, we use a tridiagonal spatial covariance matrix as an example, which can be regarded as a special case of the spherical model when sensors are located on a 1-dimensional uniform integer grid. The correlation between two sensors is ρ if they are neighboring to each other and 0 otherwise. Such a covariance matrix is denoted by $\Sigma_1(\rho) \in R_{p \times p}$ with $[\Sigma_1(\rho)]_{i,j} = 1$ if $i = j$; $[\Sigma_1(\rho)]_{i,j} = \rho$, if $|i - j| = 1$ and $[\Sigma_1(\rho)]_{i,j} = 0$, otherwise. We use $p = 5$, $\tilde{p} = 2$ and $\mu_{c,r} = [1, 1, 0, 0, 0]'$ in this example. The

ratio $m_{\text{LR}}/\tilde{m}_{\text{LR}}$ is calculated as a function of the spatial correlation ρ . As spatial correlation is unlikely to be very large in practice, we consider ρ in the range of $0 \leq \rho \leq 0.3$. A ratio smaller than one implies that the charts with full observation vectors have smaller ARL_1 than the charts with reduced observation vectors, and vice versa.

For LR charts, the ratio of ARL_1 measure for the full and reduced-dimension methods can be simplified to

$$\frac{m_{\text{LR}}}{\tilde{m}_{\text{LR}}} = \left[1 + \frac{\tilde{\boldsymbol{\mu}}'_{c,r} \boldsymbol{\Sigma}_{A^*}^{-1} \tilde{\boldsymbol{\mu}}_{c,r}}{\tilde{\boldsymbol{\mu}}'_{c,r} \boldsymbol{\Sigma}_{11}^{-1} \tilde{\boldsymbol{\mu}}_{c,r}} \right]^{-1} = \frac{2 - 6\rho^2}{2 - 5\rho^2 + 3\rho^4}.$$

We plot this ratio as a function of ρ in Figure 2(a). For LR-based charts, the ratio is always smaller than 1 as expected and it decreases as the spatial correlation ρ increases. This indicates that for LR charts, methods based on full observation vectors are always better. However, the performance loss of using reduced-dimensional vectors is small (less than 7% as shown in the plot). Thus, we expect that when the spatial correlation decays reasonably fast, in the case of known center and radius, the reduced-dimension charts do not lose much detection power compared to its full version.

For T^2 statistic based charts, there is no simple analytic expression for the ratio of ARL_1 measures. However, we may still find the ratio numerically. Figure 2(b) shows the plot of m_{T^2}/\tilde{m}_{T^2} as a function of ρ . Interestingly, the ratio is slightly greater than 1, indicating that the reduced dimension method may perform slightly better than their full dimension counterparts, for the particular covariance structure we consider in this example. However, in general, the ratio is smaller than 1 depending on the covariance structure, which we show in Section 2.4.

We also conduct simulation experiments assuming the shift cluster is known. Due to the brevity of the chapter, the results are presented in Appendix A.2. In summary, simulation results match the ARL_1 measures very well. In addition, the RD-MCUSUM charts use only 3 ~ 4% more ARL_1 than the F-MCUSUM charts for a reasonable range of spatial correlation when only a single shift cluster is considered. For T^2 charts, there is almost no performance loss for the settings we considered.

2.4 Experiments

In this section, we conduct numerical experiments under more realistic settings where the shift center and the shift radius are unknown. Then we compare the ARL_1 performance of the F- and the RD-MCUSUM charts using scan statistics.

2.4.1 Experimental setup

We consider the case in which both the shift center c and radius r are unknown. We scan over a set R of possible radii at every possible shift center. In the experiments, we use $R = \{1, \sqrt{2}\}$. In a bell-shaped signal case, we consider $R = \{0, 1, \sqrt{2}, 2, 2\sqrt{2}, 3, 3\sqrt{2}\}$.

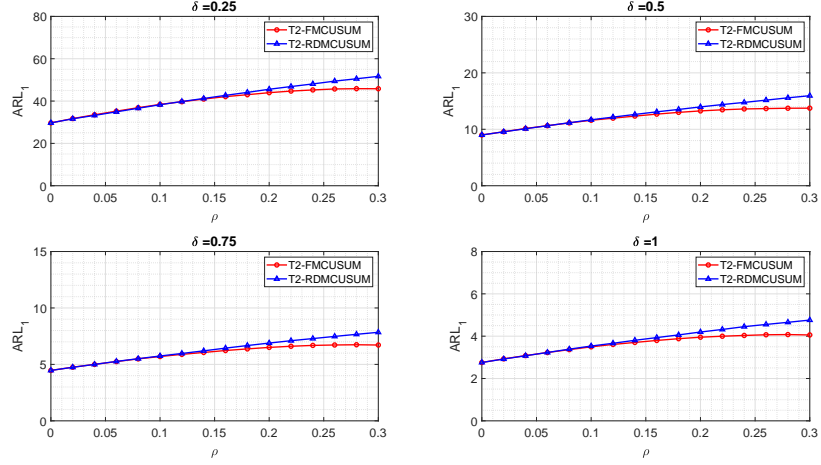
We run the control charts on three commonly used spatial models (spherical model, polynomial model and Matérn model, as introduced in Section 2.3.2). In many applications such as environmental monitoring, sensors tend to be placed with a reasonable distance and the spatial correlation coefficient between two locations is usually not high. Thus we test $\rho \in \{0, 0.02, 0.04, \dots, 0.3\}$. In the Matérn model, we test $\theta \in \{0, 0.054, 0.108, \dots, 0.81\}$ and use order $v = \frac{1}{2}$. Note that $\theta = 0.8$ for the Matérn model corresponds to the spatial correlation among neighboring regions $\rho \approx 0.3$.

The monitored area in the simulation experiments has dimensionality $p = 7 \times 7$. For the out-of-control state, homogeneous shifts (shifts of all affected locations in the cluster have same magnitude) with magnitudes $\delta = 0.25, 0.5, 0.75$ and 1 are tested. The targeted ARL_0 is fixed to 1000 in all the cases. All simulated ARL values are obtained based on 10,000 simulation replications.

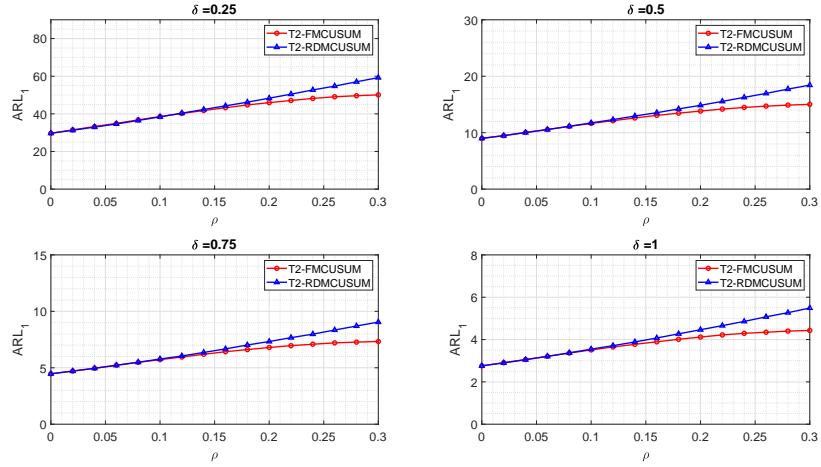
2.4.2 Results

Denote the actual radius as r_{out} . At each time step, the control chart scans over $2p$ possible shift clusters. Figures 3 compares ARL_1 of LR-F-MCUSUM and LR-RD-MCUSUM charts on three different spatial correlation structures. Figures 4 presents results for T^2 -F-MCUSUM and T^2 -RD-MCUSUM charts.

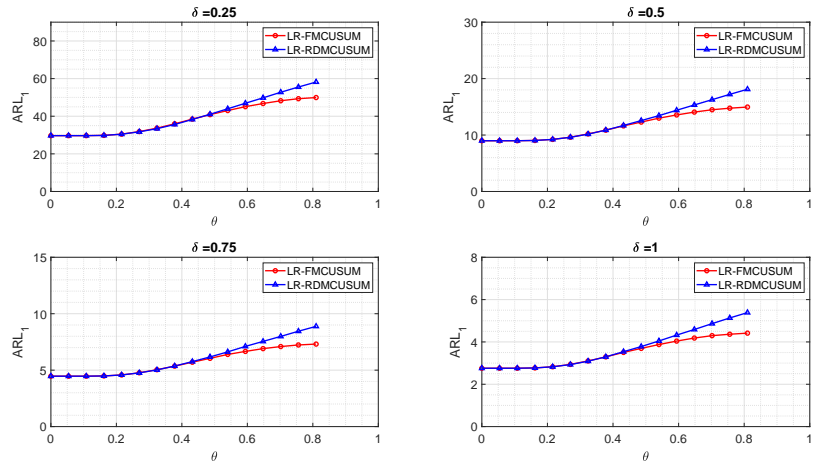
According to the numerical results, we can conclude that in general, the reduced-dimension charts do not severely sacrifice the ARL_1 performance in the range of correlation



(a)

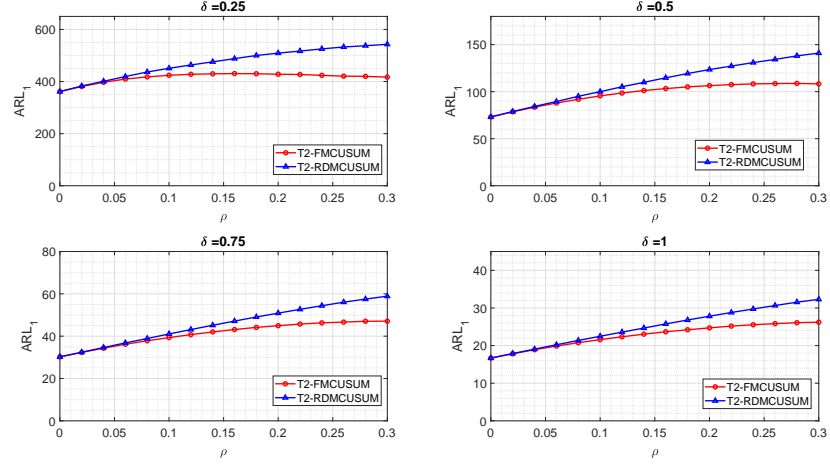


(b)

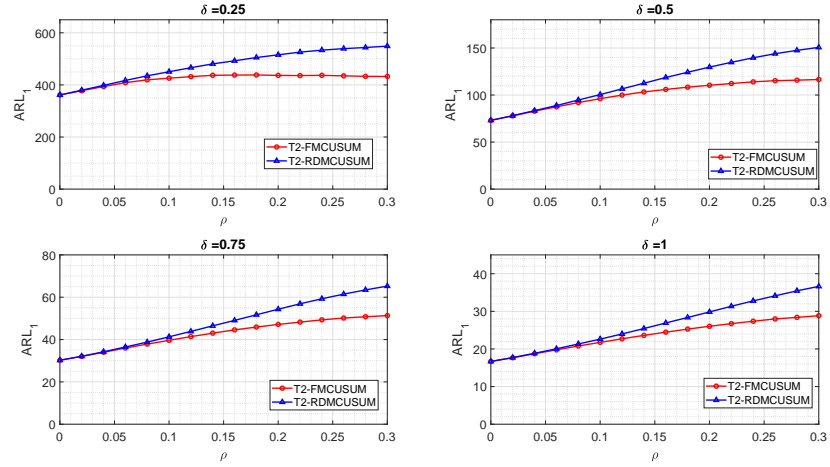


(c)

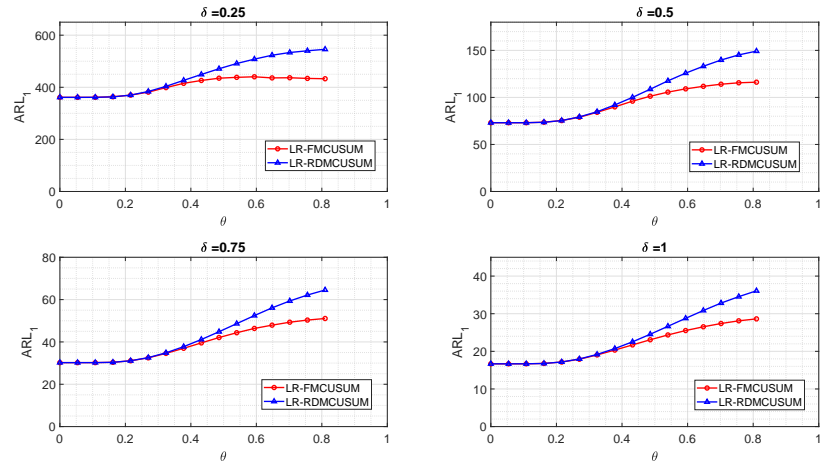
Figure 3: Simulated ARL_1 of LR-F-MCUSUM and LR-RD-MCUSUM charts with $r_{out} = \sqrt{2}$: (a) spherical model, (b) polynomial model and (c) Matérn model.



(a)



(b)



(c)

Figure 4: Simulated ARL_1 of T^2 -F-MCUSUM and T^2 -RD-MCUSUM charts with $r_{out} = \sqrt{2}$: (a) spherical model, (b) polynomial model and (c) Matérn model.

tested here. Thus, the RD-MCUSUM charts can be a powerful and easy-to-implement alternative of the F-MCUSUM charts, especially when the dimension of the monitored area is high and the full covariance matrix is ill-conditioned.

We also consider the case where the signal has a “bell” shaped rather than a boxed shape. In this example, we consider a larger set of possible shift radii, $R = \{0, 1, \sqrt{2}, 2, 2\sqrt{2}, 3, 3\sqrt{2}\}$. Appendix A.3 includes the performances of the proposed charts for this case and shows how our ARL_1 measure can be used in choosing a scanning radius in the reduced dimension charts.

2.5 Application: Water Quality Monitoring

In this section, we apply the proposed methods to real-time water quality monitoring for a river network. The goal is to detect a contaminant spill that causes water pollution in the river. In this application, the shape of the monitored region is not rectangular and scan clusters are non-circular.

2.5.1 Data

We study the Altamaha River in Georgia, United States, which is the largest watershed in the state. Figure 5 shows the shape of the Altamaha River network with 100 nodes. Each node represents a potential location to place a sensor and also a possible location of a contaminant spill. The contaminant concentration data for such a river network is simulated by the Storm Water Management Model (SWMM) developed by the United States Environmental Protection Agency. SWMM requires geologic, geometric and fundamental hydrodynamics data to construct a river network. In the simulation, rain events and spill events bring randomness to the contaminant transport. Given rainfall information, as well as the location, starting time, intensity and duration of a contaminant spill, SWMM simulates the contaminant transport process through the river over a period of time. We construct the Altamaha river system in the SWMM model based on the United States Geological Survey (USGS) digital elevation data in the National Elevation Dataset [67]. Random rain events are generated based on the patterns obtained in [65]. The Altamaha River watershed is divided into ten sub-catchments as shown in Figure 5(b). The rainfall measurements are

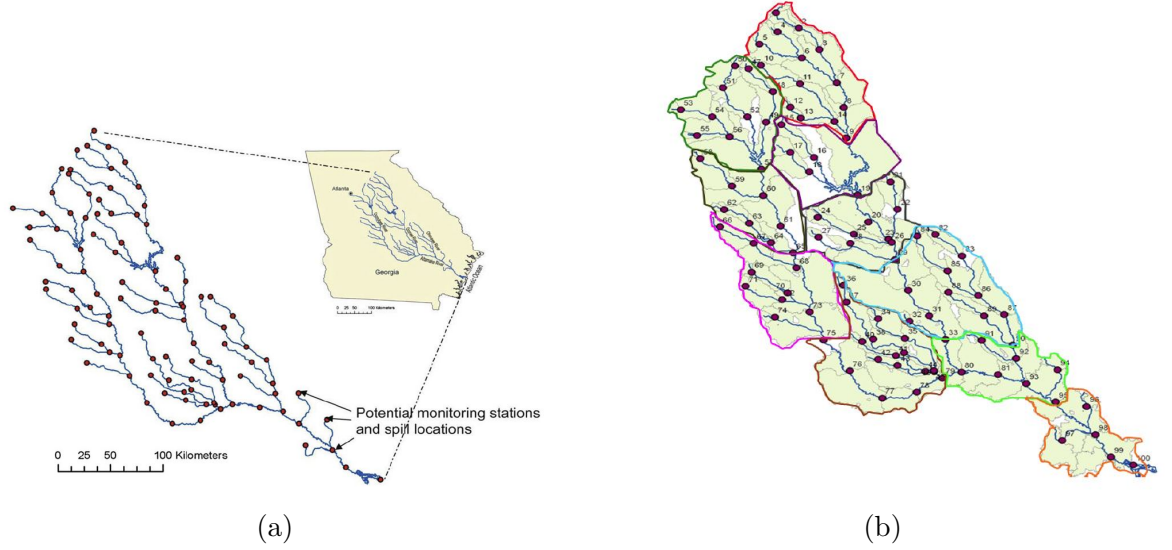


Figure 5: Shape of the Altamaha River ([65]).

obtained from different USGS stations close to the ten sub-catchments in 2006. Based on the statistical analysis of these measurements, over nine million rain patterns are generated for the entire watershed. Each rain pattern describes time-dependent rainfall events and keeps changing hydrologic conditions in each-catchment during the simulation. For each spill event, one rain pattern is randomly selected and is used in the simulation.

2.5.2 Spatial Models

Due to the nature of hydrodynamics, there exists a spatial correlation among the data streams collected at different locations in the river network. The shape of the network and direction of the stream impose constraints on modeling such spatial correlation. For example, there should not be any correlation for data collected at two monitoring locations that do not share flowing water. A reasonable spatial correlation model is critical to the detection task.

We adopt the so-called “tail-up” spatial model for stream networks, which is proposed based on the moving average constructions in [69]. The tail-up models have the following desired properties: (i) they use stream distance rather than the Euclidean distance, which is defined as the shortest distance along the stream network between two locations; (ii) statistical independence is imposed on the samples located on stream segments that do not

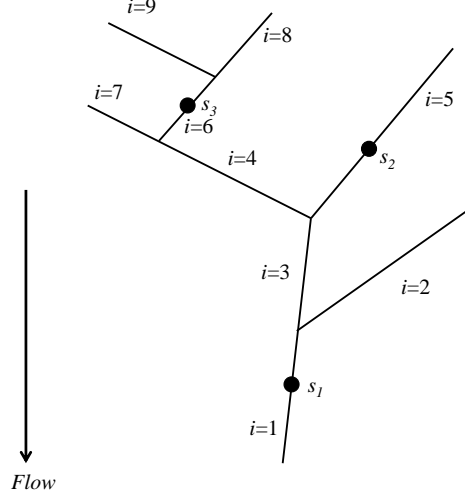


Figure 6: A stream network example with nine stream segments ($i = 1, \dots, 9$) and three monitoring locations s_1, s_2, s_3 .

share a common flow; (iii) proper weighting is incorporated on the entries of the covariance matrix when the line segments in the network is splitting into multiple segments to ensure that the resulting covariance is stationary.

To explain the tail-up model, we first introduce some notations. Suppose a stream network consists of a finite number of stream segments and we index them with $i = 1, 2, \dots$. Denote the whole set of stream segment indices as I , and the locations on the network as $s_j, j = 1, 2, \dots$. Let $D_{s_j} \subseteq I$ be the index set of all stream segments that are downstream of location s_j , (which means water from s_j flows into these segments), including the segment containing s_j . Figure 6 illustrates a simple stream network with $I = \{1, 2, \dots, 9\}$, $D_{s_1} = \{1\}$, $D_{s_2} = \{1, 3, 5\}$ and $D_{s_3} = \{1, 3, 4, 6\}$. Two locations, s_j and s_k are said to be “flow-connected” if $D_{s_j} \cap D_{s_k} = D_{s_j}$ or D_{s_k} . Finally, define

$$B_{s_j, s_k} = \begin{cases} \overline{(D_{s_j} \cap D_{s_k})} \cap (D_{s_j} \cup D_{s_k}), & \text{if } s_j \text{ and } s_k \text{ are flow-connected;} \\ \emptyset, & \text{otherwise.} \end{cases}$$

Here B_{s_j, s_k} represents the set of stream segments between two monitoring locations, including the segment for the upstream location but excluding the one for the downstream location. For example, in Figure 6, $B_{s_1, s_3} = \{3, 4, 6\}$ and $B_{s_2, s_3} = \emptyset$.

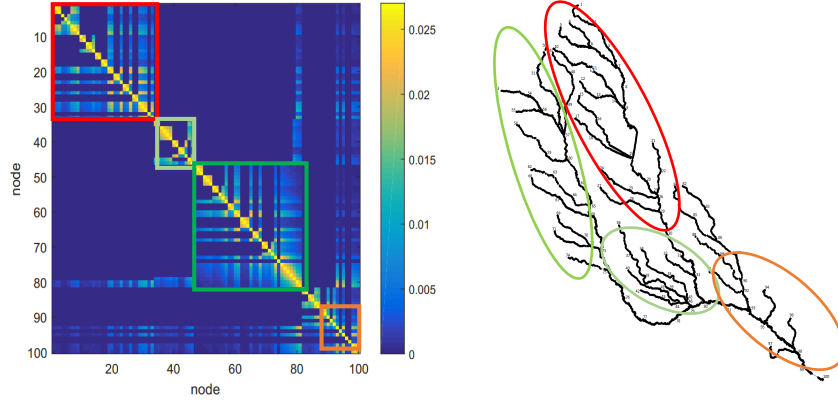


Figure 7: Visualization of the spatial covariance matrix for the Altamaha River. Each block in the covariance matrix corresponds to a branch of the river with a matching color.

To ensure the stationarity of variances, [69] suggests assigning weights to each stream segment in the network. In a stream network, one segment splits into two segments when it goes up-stream, e.g., in Figure 6, segment 1 splits into segments 2 and 3. One way to weight the segments is based on the flow volume of each segment. For example, we weight segments 2 and 3 by w_2 and w_3 , where $w_2 + w_3 = 1$ and w_2/w_3 is equal to the ratio of the flow volume between segments 2 and 3. Using tail-up models, the covariance between two locations, s_j and s_k on the stream network is given by

$$C(s_j, s_k | \zeta) = \begin{cases} 0, & \text{if } s_j \text{ and } s_k \text{ are not flow-connected;} \\ \zeta_1, & \text{if } s_j = s_k; \\ \prod_{i \in B_{s_j, s_k}} \sqrt{w_i} \zeta_1 \rho(d(s_j, s_k)/\zeta_2), & \text{otherwise;} \end{cases} \quad (14)$$

where $d(s_j, s_k)$ is the stream distance between s_j and s_k , ζ_1 is the variance parameter, $\rho(\cdot|\zeta_2)$ is the correlation function with a parameter ζ_2 , and w_i is the weight on the segment i . The correlation function $\rho(\cdot|\zeta_2)$ can be derived from many commonly used spatial models. For illustration, consider the example in Figure 6. If an exponential model is used for spatial correlation, the covariance matrix of s_1 , s_2 and s_3 can be constructed as follows:

$$\begin{pmatrix} 1 & \sqrt{w_3 w_5} & \sqrt{w_3 w_4 w_6} \\ \sqrt{w_3 w_5} & 1 & 0 \\ \sqrt{w_3 w_4 w_6} & 0 & 1 \end{pmatrix} \odot \begin{pmatrix} \zeta_1 & \zeta_1 e^{-d(s_1, s_2)/\zeta_2} & \zeta_1 e^{-d(s_1, s_3)/\zeta_2} \\ \zeta_1 e^{-d(s_1, s_2)/\zeta_2} & \zeta_1 & \zeta_1 e^{-d(s_2, s_3)/\zeta_2} \\ \zeta_1 e^{-d(s_1, s_3)/\zeta_2} & \zeta_1 e^{-d(s_2, s_3)/\zeta_2} & \zeta_1 \end{pmatrix},$$

where \odot denotes the Hadamard (element-wise) product operation between two matrices.

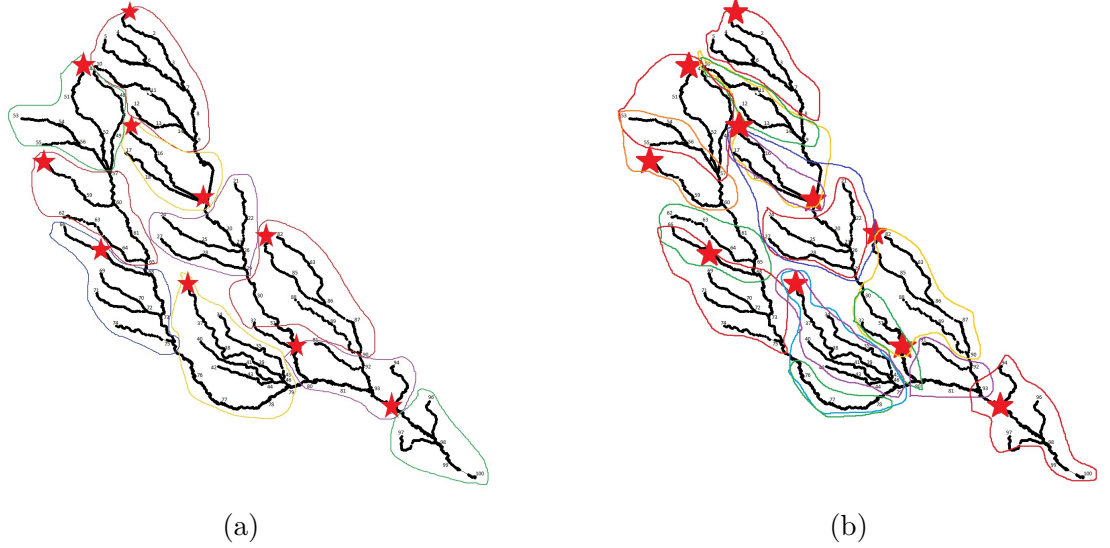


Figure 8: Two sets of scan clusters for spatial scanning: (a) non-overlapped clusters; (b) overlapped clusters. Red stars represent possible spill locations.

We use the tail-up model with an exponential correlation function to model the data collected at different nodes on the Altamaha river network. The covariance matrix for 100 nodes is constructed based on stream distances and flow volume information. We use SWMM to generate in-control data and obtain the maximum likelihood estimator of the parameters in the model, $\hat{\zeta}_1 = 0.027$ and $\hat{\zeta}_2 = 0.68$. The spatial covariance matrix is visualized in Figure 7.

2.5.3 Results

We present the detection performance of the LR-F-MCUSUM and LR-RD-MCUSUM charts for online detection of contaminant spills in the Altamaha river network. Among the 100 nodes on the river network, 10 of them are used as potential contaminant spill locations, which are marked as red stars in Figure 8, and the rest 90 nodes are used for collecting measurements every 15 minutes. In each replication, we run SWMM to simulate the river network during a 10-day period. A single instantaneous spill with a spill location randomly selected from one of the ten possible locations is generated. The spill starting time is uniformly distributed between 15 and 20 hours. Intensity of the contaminant spills follows a uniform distribution, and we consider three different levels: U(10, 100) gram/liter (low), U(100, 250) gram/liter (medium), and U(250, 500) gram/liter (high). Since the Altamaha

Table 4: Detection performance (ARL_1) obtained using the non-overlapped (NOV) and overlapped (OV) sets of scan clusters. Numbers in parentheses are standard errors.

		ARL		False alarm		Fail to detect	
	Intensity	F	RD	F	RD	F	RD
NOV	low	40.68 (4.11)	46.91 (4.60)	7%	7%	8%	9%
	medium	30.11 (2.31)	33.45 (2.90)	7%	7%	2%	9%
	high	23.79 (1.39)	25.65 (1.68)	9%	7%	0%	0%
OV	low	36.45 (3.73)	39.16 (4.22)	8%	9%	4%	8%
	medium	33.09 (2.54)	30.44 (2.26)	4%	3%	1%	3%
	high	27.30 (2.76)	27.68 (2.43)	4%	3%	1%	3%

river network does not have a regular shape and the monitoring stations are not located on the uniform grid, we do not use the circle shape clusters for spatial scanning. Instead, we construct scan clusters based on locations of the sensors and topology of the river network. Two sets of clusters are considered: (i) ten clusters that have no overlap (Figure 8(a)); and (ii) 18 clusters with partial overlap (Figure 8(b)). For the LR based MCUSUM charts, we set the minimum size of change that we aim to detect to be 0.05 gram/liter and use it to construct the post-change mean vectors, $\boldsymbol{\mu}_{c,r}$ and $\tilde{\boldsymbol{\mu}}_{c,r}$. To have a fair comparison, the thresholds for both LR-F-MCUSUM and LR-RD-MCUSUM charts are adjusted so that the in-control average run lengths are 10 days (960 samples). We generate 300 simulated contaminant spills (100 spills in each level of intensity). Detection performances of the two methods using the two sets of scan clusters are summarized in Table 4. From Table 4, we can see that the ARL performances of the two methods are consistent with our analysis in the previous sections: the chart using full observations achieves slightly smaller ARL but it is possible that the RD chart performs better as in the overlapping-medium case. To further compare the performance between the LR-F-MCUSUM and LR-RD-MCUSUM charts, we also calculate the percentage of cases where (i) the chart using full observations performs better; (ii) the two methods perform similarly, i.e., both methods successfully detect a spill event and the absolute difference of detection delays between the two methods is less than 1 hour; (iii) the chart using RD observations performs better; and (iv) both methods fail to detect by either missing the spill or raising a false alarm. These results are presented in Table 5. The table shows that, almost under all settings, the percentage

Table 5: Detection delay comparison between LR-F-MCUSUM and LR-RD-MCUSUM charts.

	non-overlapped			overlapped		
	low	medium	high	low	medium	high
Full better	35%	34%	34%	30%	18%	27%
Similar	34%	43%	45%	40%	58%	54%
RD better	25%	19%	20%	26%	22%	19%
Both fail	6%	4%	1%	4%	2%	0%

of cases where the LR-RD-MCUSUM chart performs no worse than (performs better than or similarly with) LR-F-MCUSUM chart is higher than 60%. Hence, even if the charts using full observations achieve slightly smaller ARL values, the RD based charts show very competitive performances in practice. Given that the RD based charts enjoy distributed computing with a little loss in ARL performances, they should be considered as a good alternative or even a better option for large-scale sensor networks especially when the full covariance matrix is ill-conditioned.

2.6 Conclusion

In this chapter, we study the reduced-dimension control charts for spatial-temporal change-point detection in the presence of various spatial covariance structures. The reduced-dimension charts perform spatial scanning by breaking the entire monitoring area into overlapping clusters, as well as computing the detection statistics locally while incorporating local covariance. In the presence of high-dimensional data streams, the reduced-dimension approach enjoys lower communication complexity and better numerical stability. We quantify the performance loss or gain due to dimension reduction through systematic theoretical and numerical studies. Considering the benefits of dimensionality reduction, the reduced dimension charts can be a powerful and cheaper alternative to the full-observations charts in high-dimensional problems.

CHAPTER III

S³T: A SCORE STATISTIC FOR SPATIO-TEMPORAL CHANGE-POINT DETECTION

Detection of the emergence of a signal in a noisy background arises in many multi-sensor spatio-temporal change-point detection applications. When the monitored process is in-control, sensors observe noise. When the monitored process is out of control, a signal emerges in the noise. A variety of applications possess particular spatial and temporal correlation structures. One application is an environmental sensor network, which is used to monitor river systems to detect a contaminant hazard [26]. When the signal emerges, observations from sensors may have a time-varying mean and spatio-temporal correlation structures due to water flow.

Exploiting spatio-temporal structures of the change is crucial for detecting weaker signals. However, most existing methods only capture either spatial correlations [19, 11, 22, 30, 29] or temporal correlations [73]. It is still unclear how to jointly capture the spatial and temporal correlations in detection statistics. Moreover, computational complexity is often a concern when we try to jointly model spatial and temporal correlation, especially when there are a large number of sensors that lead to high-dimensional observations. In particular, one issue with the likelihood ratio statistic is that one has to invert the sample covariance matrix, which can be computationally expensive and numerically unstable. An alternative to the likelihood ratio statistic is the score statistic, which can sometimes lead to a simpler test statistic. When the hypothesis involves a univariate parameter, the score test is the locally most powerful test [44].

In this chapter, we propose a new efficient score statistic for spatial-temporal change-point detection, which we call the S³T statistic. The S³T statistic can capture both spatial and temporal correlations of the signal. Hence, it can react quickly to a change in the mean and/or the spatio-temporal covariance. The score statistic is computationally efficient. By

avoiding the inversion of the sample covariance matrix, the $\mathbf{S}^3\mathbf{T}$ statistic has computation complexity $O(p^3)$, where p is the dimensionality of the observations, whereas the likelihood ratio statistic has $O(p^3N^3)$ complexity, which grows with the time horizon N . Our main theoretical contributions are analytic approximations for the false alarm rate in the offline case and the in-control average run length in the online case. The theoretical approximations can be used for calibrating thresholds to control the false alarm rate of our procedure. For scalar observations, our statistic $\mathbf{S}^3\mathbf{T}$ reduces to the score detector considered in [73]. Our work provides a novel extension of [73] for multi-dimensional observations when there are both spatial and temporal correlations.

The rest of the chapter is organized as follows. Section 3.1 formulates the problem. Section 3.2 presents detection statistics for both offline and online change-point detection. Section 3.3 presents theoretical approximations for the false alarm rate in the offline case and in-control average-run-length in the online case. Section 3.4 contains numerical examples for simulated data and real data, as well as a case study of water quality monitoring. Finally, Section 3.5 concludes the chapter.

3.1 Problem Formulation

Consider a sequence of samples $\mathbf{y}_\ell \in \mathbb{R}^p$, $\ell = 1, 2, \dots, N$, where p is the dimension, N is the sample size, which is fixed in the offline setting and grows in the online setting. We assume that under the null hypothesis, $\{\mathbf{y}_\ell\}$ forms a series of i.i.d. normal random vectors with spatial correlation caused by, for instance, sensor measurement errors or background noises from the environment. At an unknown time k , $1 \leq k \leq N - 1$, which corresponds to the unknown change-point, a signal emerges on top of the noise. The change may alter not only the mean of $\{\mathbf{y}_\ell\}$ but also the spatio-temporal correlation structure, which we will explain in more details.

First consider an offline setting, where the goal is to detect a change in retrospect from

the samples. Formally, this can be formulated as the following hypothesis test:

$$\begin{aligned} H_0 : & \quad \mathbf{y}_\ell = \mathbf{w}_\ell, \quad \ell = 1, 2, \dots, N, \\ H_1 : & \quad \begin{cases} \mathbf{y}_\ell = \mathbf{w}_\ell, & \ell = 1, 2, \dots, k, \\ \mathbf{y}_\ell = \mathbf{x}_\ell + \mathbf{w}_\ell, & \ell = k + 1, \dots, N, \end{cases} \end{aligned} \quad (15)$$

where $\mathbf{w}_\ell \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{\Sigma})$ and $\mathbf{\Sigma}$ is the spatial covariance matrix of the noise. We assume that before the change, the samples have no temporal correlation. This is reasonable, because we often have enough reference data before the change to estimate and then remove the temporal correlation.

Below we describe a model for the signal $\{\mathbf{x}_\ell\}$ after the change has occurred. The signal can be spatially and temporally correlated. We capture the temporal correlation using multivariate time-series models. Two examples are the first-order vector autoregressive VAR(1) model [5],

$$\mathbf{x}_\ell = (1 - \theta)\boldsymbol{\mu} + \theta\mathbf{x}_{\ell-1} + \boldsymbol{\epsilon}_\ell, \quad \ell = 1, 2, \dots, \quad (16)$$

where $\theta \in \mathbb{R}$, $\boldsymbol{\mu} = \mathbb{E}[\mathbf{x}_\ell]$ and $\boldsymbol{\epsilon}_\ell$ is the process noise; and the VARMA(1, 1) model given by

$$\mathbf{x}_{\ell+1} + \phi\mathbf{x}_\ell = (1 + \phi - \eta)\boldsymbol{\mu} + \eta\boldsymbol{\epsilon}_\ell + \boldsymbol{\epsilon}_{\ell+1}, \quad \ell = 1, 2, \dots,$$

where the parameters are $\eta \in \mathbb{R}$ and $\phi \in \mathbb{R}$. Models with higher orders can also be used if necessary.

We capture the spatial correlation of the signal using standard spatial correlation models [13]. Denote $\text{Var}[\mathbf{x}_\ell] = \gamma\mathbf{\Lambda} \in \mathbb{R}^{p \times p}$, where $\mathbf{\Lambda}$ is the spatial correlation matrix of the signal \mathbf{x}_ℓ , and $\gamma \geq 0$ is the magnitude (assuming the model is stationary and $\text{Var}[\mathbf{x}_\ell]$ does not change over ℓ). Note that the variance of the signal $\text{Var}[\mathbf{x}_\ell]$ depends on the variance of the process noise, $\text{Var}[\boldsymbol{\epsilon}_\ell]$. Here we assume the *structure* of $\mathbf{\Lambda}$ is known but the parameter γ may be unknown. This is a common practice, because once a spatial correlation model is assumed, $\mathbf{\Lambda}$ is usually specified by the location of the samples and some unknown parameters. In particular, each entry of the spatial covariance $\mathbf{\Lambda}$ is determined by a correlation function, $C(d|\rho)$, of the distance d between two samples (sensors) and is parameterized by ρ . Let $\mathbb{1}\{A\}$ denotes an indicator function, which takes value 1 when the event A is true, and 0 otherwise. Several commonly used correlation functions are:

(i) Spherical model [30]:

$$C(d|\rho) = 1\mathbb{1}\{d = 0\} + \rho\mathbb{1}\{d = 1\} + \frac{\rho}{2}\mathbb{1}\{d = \sqrt{2}\}, \quad \rho \in [0, 1]; \quad (17)$$

(ii) Exponential model [13]:

$$C(d|\rho) = 1\mathbb{1}\{d = 0\} + e^{-d/\rho}\mathbb{1}\{d > 0\}, \quad \rho > 0;$$

(iii) Matérn model [13]:

$$C(d|\rho) = 1\mathbb{1}\{d = 0\} + \frac{1}{2^{v-1}\Gamma(v)}(\sqrt{2}v^{1/2}d/\rho)^v K_v(\sqrt{2}v^{1/2}d/\rho)\mathbb{1}\{d > 0\}, \quad \rho > 0;$$

where $\rho > 0$, v is the order of the Matérn model that determines the degree of smoothness of the correlation function, $\Gamma(\cdot)$ is the gamma function, $K_v(\cdot)$ is the modified Bessel function of the second kind [46]. Note that when $v = p + 0.5$, $p \in \mathbb{R}^+$, the Matérn model is a product of an exponential and a polynomial of order p . When $v = 0.5$, the Matérn model is equivalent to the exponential model. When $v \rightarrow \infty$, it converges to the squared exponential covariance function.

Now we derive our detection statistic. For an assumed change location k , let

$$\tau = N - k$$

denote the number of post-change samples. Define a vector by concatenating all samples after the assumed change-point location k ,

$$\mathbf{y}_{(k+1:N)} = [\mathbf{y}_{k+1}^\top, \dots, \mathbf{y}_N^\top]^\top \in \mathbb{R}^{p\tau}, \quad (18)$$

where a^\top denotes the transpose of a vector a . This step is demonstrated in Figure 9. Define $\mathbf{x}_{(k+1:N)}$ and $\mathbf{w}_{(k+1:N)}$ similarly. Then after the change, we have

$$\mathbf{y}_{(k+1:N)} = \mathbf{x}_{(k+1:N)} + \mathbf{w}_{(k+1:N)}.$$

The covariance matrix of the concatenating observation vector consists of two terms that are due to the signal and the noise, respectively:

$$\text{Var}[\mathbf{y}_{(k+1:N)}] = \gamma \mathbf{V}_\tau(\theta) + \mathbf{\Sigma}_\tau,$$

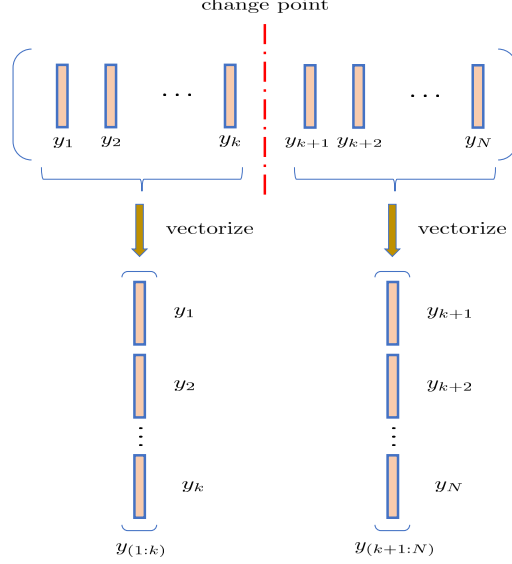


Figure 9: Diagram showing the concatenation of samples.

where $\gamma \mathbf{V}_\tau(\theta) = \text{Var}[\mathbf{x}_{(k+1:N)}]$, $\boldsymbol{\Sigma}_\tau = \text{Var}[\mathbf{w}_{(k+1:N)}]$, and θ is a parameter related to the temporal correlation which we will specify next. The second term in the covariance matrix is given by

$$\boldsymbol{\Sigma}_\tau = \mathbf{I}_\tau \otimes \boldsymbol{\Sigma} \in \mathbb{R}^{p\tau \times p\tau}, \quad (19)$$

where \mathbf{I}_τ is a τ -by- τ identity matrix and \otimes denotes the Kronecker product.

By concatenating the observation vectors as in (18), we can jointly model spatial and temporal correlation of the signal by one matrix $\mathbf{V}_\tau(\theta)$. For instance, for VAR(1) model,

$$\mathbf{V}_\tau(\theta) = \mathbf{R}_\tau(\theta) \otimes \boldsymbol{\Lambda}, \quad (20)$$

where $\mathbf{R}_\tau(\theta) \in \mathbb{R}^{\tau \times \tau}$ and $[\mathbf{R}_\tau(\theta)]_{i,j} = \theta^{|i-j|}$, $\forall i, j \in \{1, \dots, \tau\}$ is due to the temporal correlation in (16). Similarly, if the signal follows the VARMA(1,1) model, the matrix \mathbf{V} can be parameterized by $\theta \triangleq (\phi, \eta)$ with the following form:

$$\mathbf{V}_\tau(\theta) = \mathbf{R}_\tau(\phi, \eta) \otimes \boldsymbol{\Lambda}, \quad (21)$$

where $\mathbf{R}_\tau(\phi, \eta) \in \mathbb{R}^{\tau \times \tau}$; $[\mathbf{R}_\tau(\phi, \eta)]_{i,j} = 1 + \eta^2 - 2\phi\eta$, if $i = j$ and $[\mathbf{R}_\tau(\phi, \eta)]_{i,j} = \phi^{|i-j|-1}(\phi - \eta)(1 - \phi\eta)$, otherwise. For other models, similar forms of \mathbf{V}_τ can be derived: the temporal

Table 6: Notations.

p	dimension of samples
N	sample size in offline change-point detection
k	change-point location
τ	number of post change samples, $\tau = N - k$
$\mathbf{\Sigma}$	spatial covariance matrix of the noise, $\mathbf{\Sigma} = \text{Var}[\mathbf{w}_\ell]$
$\mathbf{\Lambda}$	structure of spatial covariance matrix of the signal $\text{Var}[\mathbf{x}_\ell] = \gamma \mathbf{\Lambda}$
γ	magnitude of spatial covariance matrix of the signal $\text{Var}[\mathbf{x}_\ell] = \gamma \mathbf{\Lambda}$
$\mathbf{\Sigma}_\tau$	covariance of noise in concatenated observations $\mathbf{\Sigma}_\tau = \text{Var}[\mathbf{w}_{(k+1:N)}] = \mathbf{I}_\tau \otimes \mathbf{\Sigma}$
$\gamma \mathbf{V}_\tau(\theta)$	covariance of signal in concatenated observations $\text{Var}[\mathbf{x}_{(k+1:N)}] = \gamma \mathbf{V}_\tau(\theta) = \gamma \mathbf{R}_\tau(\theta) \otimes \mathbf{\Lambda}$
$\mathbf{R}_\tau(\theta)$	matrix that captures temporal dependence of the signal

dependence of the signal is captured by \mathbf{R}_τ , the spatial dependence by $\mathbf{\Lambda}$, and the spatial-temporal covariance is a Kronecker product of the two [16].

Using the representation above, the detection problem can be reformulated as the following hypothesis test:

$$\begin{aligned} \mathbf{H}_0 : \quad & \mathbf{y}_{(1:k)} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}_k), \quad \mathbf{y}_{(k+1:N)} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}_\tau), \\ \mathbf{H}_1 : \quad & \mathbf{y}_{(1:k)} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}_k), \quad \mathbf{y}_{(k+1:N)} \sim \mathcal{N}(\boldsymbol{\mu}_{(k+1:N)}, \gamma \mathbf{V}_\tau(\theta) + \mathbf{\Sigma}_\tau), \end{aligned} \quad (22)$$

for $k = 1, \dots, N - 1$, where $\mathbf{0}$ is a vector of zeros, $\boldsymbol{\mu}_{(k+1:N)} = \mathbb{E}[\mathbf{y}_{(k+1:N)}] \in \mathbb{R}^{p\tau}$ and $\gamma \in \mathbb{R} > 0$. Note that we assume $\boldsymbol{\mu}_{(k+1:N)}$ is unknown. Equivalently, under the null hypothesis, $\gamma = 0$ and $\boldsymbol{\mu}_{(k+1:N)} = \mathbf{0}$, and under the alternative hypothesis, $\gamma > 0$ or $\boldsymbol{\mu}_{(k+1:N)} \neq \mathbf{0}$. Using this form of hypothesis, we can derive our score statistic.

Table 6 provides a list of notations used throughout the chapter.

3.2 Statistic for Offline and Online Detection

We now derive the $\mathbf{S}^3\mathbf{T}$ statistic for offline change-point detection. The log-likelihood function for the hypothesis test in (22) is given by

$$\begin{aligned} \ell(\gamma, \boldsymbol{\mu}, \tau, \theta) = & -\frac{1}{2} \log(2\pi) - \frac{1}{2} \log |\gamma \mathbf{V}_\tau(\theta) + \mathbf{\Sigma}_\tau| \\ & - \frac{1}{2} (\mathbf{y}_{(k+1:N)} - \boldsymbol{\mu}_{(k+1:N)})^\top (\gamma \mathbf{V}_\tau(\theta) + \mathbf{\Sigma}_\tau)^{-1} (\mathbf{y}_{(k+1:N)} - \boldsymbol{\mu}_{(k+1:N)}). \end{aligned} \quad (23)$$

To cope with unknown parameters, we may use the generalized likelihood ratio (GLR) statistic based on (23). However, (23) involves the inversion of a $p\tau$ -by- $p\tau$ dimensional matrix $\gamma \mathbf{V}_\tau(\theta) + \mathbf{\Sigma}_\tau$, which incurs a complexity of $O(p^3\tau^3)$ for a given τ . Recall that

$\tau = N - k$, so $\tau = 1, 2, \dots, N$. Hence, the complexity of computing the GLR statistic is $O(p^3 N^3)$, which grows polynomially with N (the time horizon), and the computation of the likelihood statistic becomes prohibitive.

3.2.1 Quadratic score statistic

Define the following notations. Let $\mathbf{A}_\tau(\theta) = \mathbf{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta)$, $\mathbf{B}_\tau(\theta) = \mathbf{\Sigma}_\tau^{-1/2} \mathbf{V}_\tau(\theta) \mathbf{\Sigma}_\tau^{-1/2}$, $c(\tau, \theta) = \text{tr}(\mathbf{A}_\tau(\theta))$, and $d(\tau, \theta) = 2\text{tr}[\mathbf{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta) \mathbf{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta)]$, where $\text{tr}(\cdot)$ denotes the trace of a matrix. We now derive the score-statistic for detection. Take the derivative of $\ell(\gamma, \boldsymbol{\mu}, \tau, \theta)$ in (23) with respect to γ and $\boldsymbol{\mu}$ and evaluate at $\gamma = 0$ and $\boldsymbol{\mu} = \mathbf{0}$. We can obtain

$$\varsigma(\tau, \theta) = \begin{bmatrix} \frac{\partial \ell}{\partial \gamma} \big|_{\boldsymbol{\mu}=\mathbf{0}, \gamma=0} \\ \frac{\partial \ell}{\partial \boldsymbol{\mu}} \big|_{\boldsymbol{\mu}=\mathbf{0}, \gamma=0} \end{bmatrix} = \begin{bmatrix} -\frac{1}{2}c(\tau, \theta) + \frac{1}{2} \mathbf{y}_{(k+1:N)}^\top \mathbf{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta) \mathbf{\Sigma}_\tau^{-1} \mathbf{y}_{(k+1:N)} \\ \mathbf{\Sigma}_\tau^{-1} \mathbf{y}_{(k+1:N)} \end{bmatrix}. \quad (24)$$

The derivation of (24) is given in Appendix B.1. It can be verified that $\mathbb{E}[\varsigma(k, \theta)] = \mathbf{0}$ under the null hypothesis, where $\mathbf{0}$ represents the zero vector. It can also be shown that the covariance of the score vector $\varsigma(\tau, \theta)$ is given by

$$\text{Cov}[\varsigma(\tau, \theta)] = \begin{bmatrix} \frac{1}{4}d(\tau, \theta) & 0 \\ \mathbf{0} & \mathbf{\Sigma}_\tau^{-1} \end{bmatrix}.$$

As suggested by [43], when the likelihood function involves multiple parameters, the score statistic is a quadratic function of the efficient score vector. In our case, this becomes

$$\begin{aligned} S(\tau, \theta) &= \varsigma(\tau, \theta)^\top \text{Cov}[\varsigma(\tau, \theta)]^{-1} \varsigma(\tau, \theta) \\ &= \frac{\left[\mathbf{y}_{(k+1:N)}^\top \mathbf{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta) \mathbf{\Sigma}_\tau^{-1} \mathbf{y}_{(k+1:N)} - c(\tau, \theta) \right]^2}{d(\tau, \theta)} + \mathbf{y}_{(k+1:N)}^\top \mathbf{\Sigma}_\tau^{-1} \mathbf{y}_{(k+1:N)}. \end{aligned} \quad (25)$$

The most expensive part in evaluating (25) is to compute $\mathbf{\Sigma}_\tau^{-1}$. According to (19), we have $\mathbf{\Sigma}_\tau^{-1} = \mathbf{I}_\tau \otimes \mathbf{\Sigma}^{-1}$, which means that we only need to compute $\mathbf{\Sigma}^{-1}$ that has a complexity $O(p^3)$. Hence, the computation complexity of evaluating $S(\tau, \theta)$ is much lower than that of the GLR statistic. Moreover, since $\mathbf{\Sigma}$ is assumed known and fixed, its inversion can be pre-computed; however, in (23), the likelihood function involves $(\gamma \mathbf{V}_\tau(\theta) + \mathbf{\Sigma}_\tau)^{-1}$, which has to be computed for each τ value.

Since the expected value of $S(\tau, \theta)$ increases as τ increases, it needs to be normalized to have mean 0 and variance 1 under the null hypothesis. This leads to the following *quadratic*

score statistic,

$$\tilde{S}(\tau, \theta) = \frac{S(\tau, \theta) - \mathbb{E}[S(\tau, \theta)]}{\sqrt{\text{Var}[S(\tau, \theta)]}}, \quad (26)$$

where $\mathbb{E}[S(\tau, \theta)] = p\tau + 1$, and the variance is given by (derivation can be found in Appendix B.2)

$$\begin{aligned} \text{Var}[S(\tau, \theta)] = & 2p\tau + 10 - 24 \frac{c(\tau, \theta)}{d(\tau, \theta)^2} \text{tr} \left(\mathbf{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta) \mathbf{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta) \mathbf{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta) \right) \\ & + \frac{48}{d(\tau, \theta)^2} \text{tr} \left(\mathbf{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta) \mathbf{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta) \mathbf{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta) \mathbf{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta) \right). \end{aligned}$$

Then we may construct the *quadratic detector* using $\tilde{S}(\tau, \theta)$, which detects a signal when the maximum standardized score statistic over all possible parameter values of $\theta \in \Theta$ and τ exceeds a pre-specified threshold $b > 0$,

$$\max_{\theta \in \Theta, 1 \leq \tau \leq N} \tilde{S}(\tau, \theta) \geq b.$$

3.2.2 S^3T statistic for offline change-point detection

Although the quadratic score statistic achieves the maximum discrimination between the null and the alternative distribution [43], theoretical analysis of the detection statistic is intractable; thus, it is difficult to calibrate the threshold b to control the false alarm rate. In this section, we propose a simpler statistic, namely the S^3T statistic, which is the score statistic with respect to γ *only*:

$$W(\tau, \theta) = \frac{\frac{\partial \ell}{\partial \gamma} \big|_{\boldsymbol{\mu}=\mathbf{0}, \gamma=0}}{\sqrt{\text{Var}[\frac{\partial \ell}{\partial \gamma} \big|_{\boldsymbol{\mu}=\mathbf{0}, \gamma=0}]}} = \frac{\mathbf{y}_{(N-\tau+1:N)}^\top \mathbf{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta) \mathbf{\Sigma}_\tau^{-1} \mathbf{y}_{(N-\tau+1:N)} - c(\tau, \theta)}{\sqrt{d(\tau, \theta)}}. \quad (27)$$

Note that both the spatial and temporal correlations of the signal are still captured in the statistic by $\mathbf{V}_\tau(\theta)$ in (20). Under the null hypothesis, the detection statistic $W(\tau, \theta)$ has mean 0 and unit variance. The detection procedure claims a change when the maximum of the score statistic exceeds a pre-specified threshold $b > 0$,

$$\max_{\theta \in \Theta, 1 \leq \tau \leq N} W(\tau, \theta) \geq b. \quad (28)$$

3.2.3 S^3T statistic for online change-point detection

We now present an online change-point detection procedure based on the S^3T statistic. In the online setting, the sample size N is not fixed and samples are sequentially collected. A

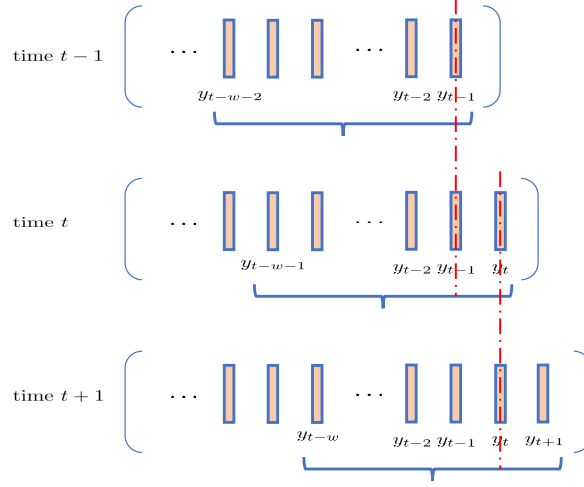


Figure 10: Sliding window of length w for online detection.

signal may occur at an unknown time t . Our goal is to detect the emergence of the signal as soon as possible.

Consider a sequential version of the hypothesis test in (15), where the number of samples N increase. We adopt a sliding window approach for online detection, and construct the detection statistic using the most recent ω samples at each time, where ω is a pre-specified window length (demonstrated in Figure 10). Given a current time t , the detection statistic constructed using the most recent ω samples is given by

$$W_t(\omega, \theta) = \frac{\mathbf{y}_{(t-\omega+1:t)}^\top \boldsymbol{\Sigma}_\omega^{-1} \mathbf{V}_\omega(\theta) \boldsymbol{\Sigma}_\omega^{-1} \mathbf{y}_{(t-\omega+1:t)} - c(\omega, \theta)}{\sqrt{d(\omega, \theta)}}. \quad (29)$$

The detection procedure for online change-point detection is a stopping time, which raises an alarm when the detection statistic exceeds a threshold $b > 0$ for the first time:

$$\mathcal{T} = \inf \left\{ t : \max_{\theta \in \Theta} W_t(\omega, \theta) \geq b \right\}. \quad (30)$$

3.3 Theoretical Approximations

3.3.1 Significance level for offline S^3T statistic

We present a theoretical approximation for the significance level of the detection procedure defined in (28). The approximation is quite accurate and can be used to avoid the time-consuming simulation when choosing an appropriate b . Denote the standard normal density

function by $\phi(x)$ and its distribution function by $\Phi(x)$, and define a special function [53]:

$$\nu(x) = \frac{\frac{2}{x} \left[\Phi\left(\frac{x}{2}\right) - \frac{1}{2} \right]}{\frac{x}{2} \Phi\left(\frac{x}{2}\right) + \phi\left(\frac{x}{2}\right)}. \quad (31)$$

Define the following quantities, which are useful to state our theoretical approximation results:

$$\mu(\tau, \theta) = \tau \left[\frac{\text{tr}(\mathbf{A}_{\tau+1}(\theta) \mathbf{A}_{\tau+1}(\theta))}{\text{tr}(\mathbf{A}_\tau(\theta) \mathbf{A}_\tau(\theta))} - 1 \right], \quad (32)$$

$$H(\tau, \theta) = - \frac{\partial^2 E[W(\tau, \theta) W(\tau, s)]}{\partial^2 s} \Big|_{s=\theta}, \quad (33)$$

$$\psi(\xi) = -\xi \frac{c(\tau, \theta)}{\sqrt{d(\tau, \theta)}} - \frac{1}{2} \log \left| \mathbf{I}_{p\tau} - \frac{2\xi \mathbf{B}_\tau(\theta)}{\sqrt{d(\tau, \theta)}} \right|. \quad (34)$$

Note that $\psi(\xi)$ is the cumulant generating function (a.k.a., the log-moment generating function) of the detection statistic $W(\tau, \theta)$. The following theorem is our main theoretical result, which provides an analytical approximation for the significance level of the detection procedure defined in (28).

Theorem 3.3.1 (Approximation for significance level). *When the threshold $b \rightarrow \infty$ and $\theta \in \Theta \subset \mathbb{R}^d$, under the null hypothesis, the probability of false alarm for the procedure defined in (28) is given by*

$$\begin{aligned} & \mathbb{P}_{H_0} \left(\max_{\substack{\theta \in \Theta \\ 1 \leq \tau \leq N}} W(\tau, \theta) \geq b \right) \\ &= \frac{1}{(2\pi)^{\frac{d}{2}}} \sum_{\tau=1}^N \int_{\theta \in \Theta} \frac{[b\xi_0(\tau, \theta)]^{\frac{d}{2}}}{\xi_0(\tau, \theta)} g(\tau, \theta) |H(\tau, \theta)|^{\frac{1}{2}} \frac{b^2 \mu(\tau, \theta)}{2\tau} \nu\left(\sqrt{\frac{b^2 \mu(\tau, \theta)}{\tau}}\right) d\theta + o(1), \end{aligned} \quad (35)$$

where

$$g(\tau, \theta) = \frac{\exp(-\xi_0(\tau, \theta)b + \psi(\xi_0(\tau, \theta)))}{\sigma_{\xi_0} \sqrt{2\pi}}, \quad (36)$$

$$\sigma_{\xi_0}^2 = d(\tau, \theta)^{-1} \text{tr} \left(\left[\mathbf{I}_{p\tau} - \frac{2\xi_0 \mathbf{B}_\tau(\theta)}{\sqrt{d(\tau, \theta)}} \right]^{-1} \mathbf{B}_\tau(\theta) \left[\mathbf{I}_{p\tau} - \frac{2\xi_0 \mathbf{B}_\tau(\theta)}{\sqrt{d(\tau, \theta)}} \right]^{-1} \mathbf{B}_\tau(\theta) \right), \quad (37)$$

and $\xi_0(\tau, \theta)$ is the solution to

$$\frac{1}{\sqrt{d(\tau, \theta)}} \text{tr} \left(\left[\mathbf{I}_{p\tau} - \frac{2\xi_0 \mathbf{B}_\tau(\theta)}{\sqrt{d(\tau, \theta)}} \right]^{-1} \mathbf{B}_\tau(\theta) - \mathbf{A}_\tau(\theta) \right) = b. \quad (38)$$

Table 7: Simulated and approximated significance level when the signal $\{\mathbf{x}_\ell\}$ follows a VAR(1) model.

b	$p = 2$		$p = 9$		$p = 36$	
	Simulated	Approximated	Simulated	Approximated	Simulated	Approximated
3.5	0.097	0.097	0.065	0.057	0.036	0.042
4	0.063	0.068	0.036	0.030	0.013	0.019
4.5	0.038	0.047	0.018	0.019	0.006	0.008
5	0.033	0.032	0.011	0.012	0.003	0.003
5.5	0.022	0.021	0.005	0.007	0.002	0.001
6	0.015	0.014	0.003	0.004	0.0004	0.0005
6.5	0.006	0.009	0.002	0.002	0.0002	0.0002

Note that the solution of (38) can be obtained by a simple grid search when the dimension of θ is not too large.

The main proof technique for Theorem 3.3.1 is *change-of-measure*, which evaluates the boundary hitting probability of Gaussian processes [55, 76]. See Appendix B.3 for the derivation of (34) and Appendix B.4 for the proof of Theorem 3.3.1, when the dimension of parameter θ is 1 (i.e., $d = 1$). The proof can be generalized to multi-dimensional cases.

Although the theorem is an asymptotic result for large b , we find that this holds even for not very large b values in numerical studies. We verify the accuracy of Theorem 3.3.1 by comparing the approximated significance levels with simulated ones. In the experiment, we assume that the signal $\{\mathbf{x}_\ell\}$ follows a VAR(1) model, $\mathbf{x}_\ell = (1 - \theta)\boldsymbol{\mu} + \theta\mathbf{x}_{\ell-1} + \boldsymbol{\epsilon}_\ell$, where $\theta \in \mathbb{R}$. Hence, $\mathbf{V}_\tau(\theta)$ has the form in (20). We further assume that the spatial correlation of the signal follows a spherical model, as defined in (17), with parameter $\rho = 0.3$. Set $N = 50$. The search space of θ is $\{0.1, 0.2, \dots, 0.9\}$. In addition, the covariance matrix of the noise process $\boldsymbol{\Sigma}$ is assumed to be a p -by- p identity matrix. Simulation results are based on 5000 independent replications. Both simulated and approximated false alarm rates are reported in Table 7. As one can observe, the approximation is quite accurate.

In the proof of Theorem 3.3.1, we approximate the detection statistic $W(\tau, \theta)$ as a two-dimensional Gaussian random field. In the following, we verify that such an approximation is reasonable by simulation. We generate data under the null hypothesis, and verify the distribution of the detection statistic W for a set of fixed values of θ and τ . For the signal, we use a VAR(1) model, $\mathbf{x}_\ell = (1 - \theta)\boldsymbol{\mu} + \theta\mathbf{x}_{\ell-1} + \boldsymbol{\epsilon}_\ell$ as the temporal correlation model and

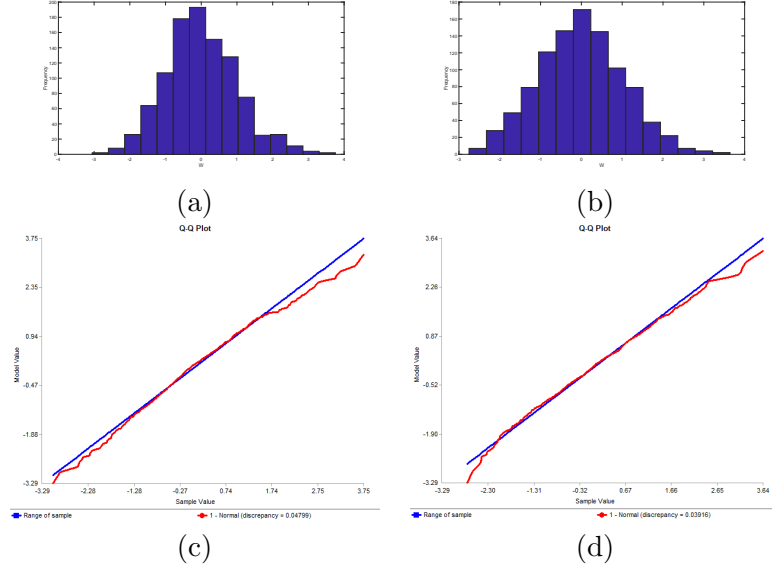


Figure 11: Histograms and q-q plots of $W(\theta, \tau)$ for fixed values of θ and τ : $\tau = 30$, $\theta = 0.3$ for (a) and (c); $\tau = 40$, $\theta = 0.2$ for (b) and (d).

a spherical model for spatial correlation model. We assume that the noise has the same spatial correlation structure as the signal. We set $N = 50$ and $p = 9$. Figure 11 shows the histograms and q-q plots of W for fixed values of θ and τ based on 1000 repetitions, which indicate that the Gaussian random field approximation is reasonable.

3.3.2 In-control Average Run Length (ARL_0) for online S^3T statistic

In the online setting, the false alarm rate is characterized by the in-control average-run-length, which is equal to the expected stopping time of the procedure when there is no signal, denoted as $E_{H_0}(\mathcal{T})$. The following theorem provides an approximation for $E_{H_0}(\mathcal{T})$.

Theorem 3.3.2 (Approximation of ARL_0). *Assume that $b \rightarrow \infty$. For the stopping time defined in (30),*

$$E_{H_0}(\mathcal{T}) = (2\pi)^{\frac{d}{2}} \left(\int_{\theta \in \Theta} \frac{[b\xi_0(\omega, \theta)]^{\frac{d}{2}}}{\xi_0(\omega, \theta)} g(\omega, \theta) |H(\omega, \theta)|^{\frac{1}{2}} \frac{b^2 \mu(\omega, \theta)}{2\omega} \nu\left(\sqrt{\frac{b^2 \mu(\omega, \theta)}{\omega}}\right) d\theta \right)^{-1} (1 + o(1)). \quad (39)$$

The derivation of Theorem 3.3.2 uses a similar technique based on the change-of-measure

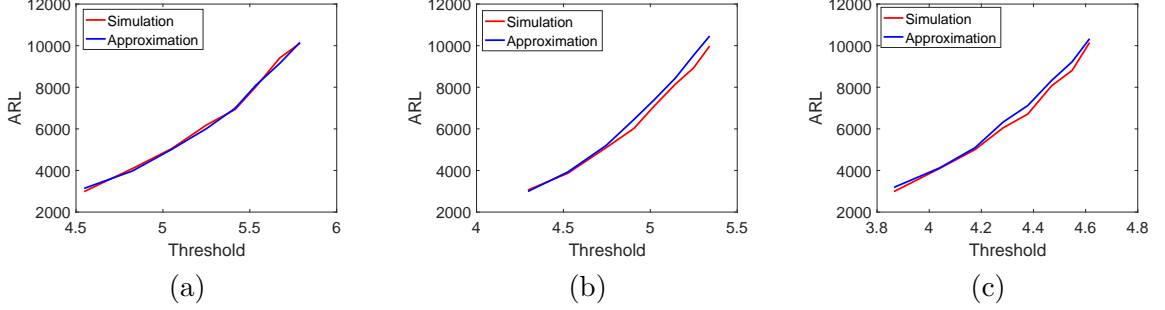


Figure 12: Comparison of approximated and simulated ARL for (a) $p = 1$, (b) $p = 2$, and (c) $p = 9$.

as in the derivation of Theorem 3.3.1. By Theorem 3.3.1, we can first obtain an approximation for the probability $\mathbb{P}_{H_0}(\mathcal{T} \leq m)$, where $m > 0$ is fixed and sufficiently large:

$$\begin{aligned} \mathbb{P}_{H_0}(\mathcal{T} \leq m) &= \mathbb{P}_{H_0}\left(\max_{\substack{\theta \in \Theta \\ 1 \leq t \leq m}} W_t(\omega, \theta) \geq b\right) \\ &= (2\pi)^{-\frac{d}{2}} \left(\sum_{t=1}^m \int_{\theta \in \Theta} \frac{[b\xi_0(\omega, \theta)]^{\frac{d}{2}}}{\xi_0(\omega, \theta)} g(\omega, \theta) |H(\omega, \theta)|^{\frac{1}{2}} \frac{b^2 \mu(\omega, \theta)}{2\omega} \nu\left(\sqrt{\frac{b^2 \mu(\omega, \theta)}{\omega}}\right) d\theta \right) + o(1). \end{aligned} \quad (40)$$

As argued in [56] and [57], the stopping time \mathcal{T} is asymptotically exponentially distributed and is uniformly integrable. Hence, for large m , $\mathbb{P}_{H_0}(\mathcal{T} \leq m) - [1 - \exp(-\lambda m)] \rightarrow 0$, where λ is approximately equal to the right hand side of (40) divided by m . Then by the first order Taylor expansion of an exponential term, we can obtain $E_{H_0}(\mathcal{T}) \approx \lambda^{-1}$, which leads to (39).

The accuracy of Theorem 3.3.2 is verified by comparing the simulated and the approximated $E_{H_0}(\mathcal{T})$. In the experiments, the signal $\{\mathbf{x}_\ell\}$ is generated by a VAR(1) model, $\mathbf{x}_\ell = (1 - \theta)\boldsymbol{\mu} + \theta\mathbf{x}_{\ell-1} + \boldsymbol{\epsilon}_\ell$, where $\theta \in \mathbb{R}$. Hence, $\mathbf{V}_\tau(\theta)$ has the form in (20). Meanwhile, we assume that the spatial correlation of the signal follows a spherical model, as defined in (17), with parameter $\rho = 0.3$. The search space of parameter θ is $\{0.1, 0.2, \dots, 0.9\}$. In addition, the covariance matrix of the noise process $\boldsymbol{\Sigma}$ is assumed to be a p -by- p identity matrix. The results based on 5000 replications are presented in Figure 12. The comparison between the simulated and approximated ARLs shows that the approximation in Theorem 3.3.2 is quite accurate.

3.4 Numerical Examples

In this section, we demonstrate the performance of the proposed detection procedures. Online change-point detection is the focus here, since it is the most relevant setting for our targeted applications of water quality monitoring. The performance comparison for offline change-point detection will be similar. We adopt the commonly used performance metric for sequential change detection, the expected detection delay (EDD) after a change has occurred. There is a tradeoff between the in-control average-run-length (ARL_0) and the EDD. Typically, we choose the threshold for each procedure so that its ARL_0 meets a pre-specified large value (e.g., 5000 or 10000), and hence there is rarely a false alarm. We also compare with other methods on simulated and real data.

3.4.1 Simulation

The detection procedure defined in (30) is compared with two other procedures: (i) an online detection procedure defined similarly to (30) using the quadratic score statistic $\tilde{S}(\tau, \theta)$, and (ii) a multivariate cumulative sum (MCUSUM) procedure [19]. In the MCUSUM procedure, at each time step, a T^2 statistic [21] is calculated, which is combined with a CUSUM procedure.

In the experiment, the signal is generated from a VAR(1) model, $\mathbf{x}_\ell = (1 - \theta)\boldsymbol{\mu} + \theta\mathbf{x}_{\ell-1} + \boldsymbol{\epsilon}_\ell$, with $p = 2$ and parameter $\theta = 0.5$. The spatial model of the signal follows the spherical model defined in (17) with $\rho = 0.3$. For both procedures, based on $\mathbf{S}^3\mathbf{T}$ and the quadratic score statistic, respectively, we use a window length $\omega = 50$ and the search space for the parameter θ , $\{0.1, 0.2, \dots, 0.9\}$. Thresholds for all three procedures are calibrated so that they have the same false alarm rate $E_{H_0}(\mathcal{T}) = 100$. To evaluate the expected detection delay, we assume the change occurs at $t = 1$. The mean of the signal $\boldsymbol{\mu} = E[\mathbf{x}_\ell] = \mu \mathbb{1}_p$, $\mu \geq 0$. We explore different values of μ for the mean shift and γ for the magnitude of covariance matrix of the signal. If $\mu = 0$ and $\gamma > 0$, there is only a change in covariance; if both μ and γ are positive, then there are both mean shift and covariance change. Hence, the experiments demonstrate that the proposed detection procedure is suitable for both cases where there is either mean and/or covariance change.

Table 8: Simulated expected detection delay.

	S^3T					Quadratic score statistic					MCUSUM				
$\gamma \backslash \mu$	0	0.1	0.5	1	2	0	0.1	0.5	1	2	0	0.1	0.5	1	2
0.01	97.27	59.08	6.37	2.80	1.49	98.05	65.82	6.45	2.77	1.51	98.37	77.67	9.43	3.56	1.79
0.05	96.28	57.96	5.95	2.72	1.49	95.32	63.19	6.74	2.81	1.52	96.79	71.97	9.28	3.54	1.79
0.1	72.93	53.16	6.04	2.78	1.50	82.49	56.78	6.74	2.86	1.49	80.70	65.16	9.21	3.54	1.78
0.2	65.32	46.16	5.96	2.77	1.50	74.87	48.83	6.28	2.78	1.47	67.33	55.17	9.02	3.52	1.79
0.5	39.40	30.32	5.81	2.78	1.56	37.07	33.42	6.07	2.80	1.50	41.52	35.87	8.36	3.47	1.78
1	20.91	19.42	5.65	2.75	1.51	22.75	20.51	5.64	2.76	1.55	23.71	21.31	7.45	3.45	1.77

Table 8 reports the simulated EDD of three procedures based on 5000 repetitions. The smallest EDD values for each setting are marked in bold. The comparison shows that the two score statistic procedures, which capture both spatial and temporal correlation, outperform the MCUSUM procedure (which only captures the spatial correlation information). Such an advantage is more significant when the signal is weak, i.e., when γ or μ are both small. This demonstrates that incorporating temporal correlation information indeed improves detection performance. We also find that S^3T outperforms the quadratic score statistic in many settings. This can be explained by that the quadratic score statistic needs to search more unknown parameters (the unknown μ), thus the statistic is noisier than S^3T when there is no change. Therefore, to achieve the same ARL_0 , the threshold for quadratic score statistic tends to be higher, which may cause a larger detection delay. Given that S^3T enjoys tractable theoretical analysis and an accurate approximation for its false alarm rate, it is a good option for practitioners.

3.4.2 Real data example: Solar flare detection

We apply our detection procedure to a real dataset, which is acquired by the Solar Data Observatory [37]. The data is a video sequence that contains an abrupt emergence of a solar flare that occurs around time $t = 227$. In this video, the normal state is a sequence of slowly drifting image of the solar surface, and the changes are much brighter transient solar flares. Figure 13 shows a snapshot when a solar flare occurs at $t = 227$.

The size of the images is 232×292 pixels. After vectoring the images, this leads to 67,744 dimensional vectors. Due to the high dimensionality, it is computationally expensive to apply our detection procedure on the original images directly. Hence, we apply the

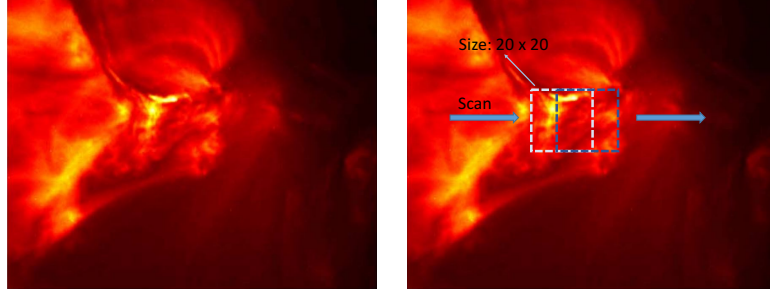


Figure 13: Detection of solar flare at $t = 227$: (left) snapshot of the original SDO data at $t = 227$; (right) overlapping image patches for dimensionality reduction.

spatial scanning scheme proposed in the previous chapter. We break the original image into overlapping patches of dimension 20×20 , as demonstrated in the right panel of Figure 13. The detection statistic is calculated for each image patch (of dimension $p = 400$). Then, we take the maximum of the detection statistic over all patches.

We assume that before the solar flare, the data form a white noise process with no spatial and temporal correlation. The mean and variance of the noise process are estimated by the first 50 samples in the sequence. For the signal, we use a VAR(1) model, $\mathbf{x}_\ell = (1 - \theta)\boldsymbol{\mu} + \theta\mathbf{x}_{\ell-1} + \boldsymbol{\epsilon}_\ell$ to capture the temporal correlation. The spatial model of the signal is captured by a spherical model defined in (17). Online procedures are implemented with window length $\omega = 10$. Figure 14(a), 14(b) and 14(c) show the values of $\mathbf{S}^3\mathbf{T}$ statistic, the quadratic score statistic and MCUSUM statistic on a logarithmic scale, respectively. Since in this case, we do not have the ground truth, we cannot evaluate the true EDD. However, as we can observe, both $\mathbf{S}^3\mathbf{T}$ and the quadratic score statistics obtain peak detection statistics at around $t = 227$, and another solar flare at around $t = 173$, indicating both statistics can successfully detect the emergence of solar flares. However, MCUSUM statistic misses both solar flares.

3.4.3 Case study: Water quality monitoring

In this section, we consider a case study of real-time water quality monitoring for the Altamaha River network based on data generated by the SWMM model. The goal is to detect contaminant spills that pollute the river as quickly as possible. Backgrounds on the

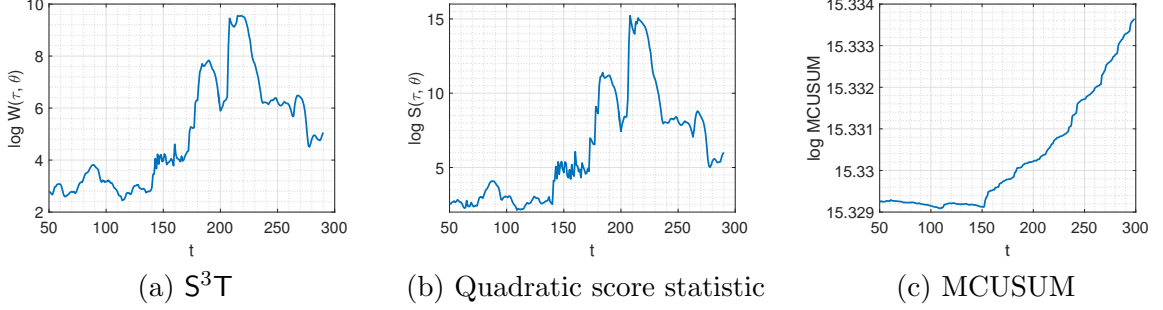


Figure 14: Detection statistics on logarithmic scale.

Altamaha River and SWMM model are presented in Section 2.5.1.

In the case study, the tail-up model (discussed in Section 2.5.2) with an exponential correlation function is adopted as the spatial model for the data. We assume that both the signal and the noise share the same spatial correlation structure. For temporal correlation, we use a VAR(1) model $\mathbf{x}_\ell = (1 - \theta)\boldsymbol{\mu} + \theta\mathbf{x}_{\ell-1} + \boldsymbol{\epsilon}_\ell$ to capture the temporal correlation of a contaminant spill as suggested in [8] and [9]. We apply the online change-point detection procedure based on S^3T to detect contaminant spills in the Altamaha River network. We also compare it with two other methods: (i) online detection based on the quadratic score statistic, and (ii) the Hotelling's T^2 chart. Among the 100 nodes on the river network, 10 of them (nodes 1, 15, 19, 33, 36, 50, 58, 67, 84, 95, marked by red stars in Figure 8(c)) are used as possible contaminant spill locations, and the rest 90 nodes are used for collecting measurements every 15 minutes. In each replication, we run SWMM to simulate the river network during a 10-day period. A single instantaneous spill is generated, with a spill location randomly selected from the ten possible locations. The spill starting time is uniformly distributed between the first 15 to 20 hours. The intensity of the contaminant spills follows a uniform distribution, and we consider three different levels: $U(10, 100)$ (low), $U(100, 250)$ (medium), and $U(250, 500)$ (high) in units of gram/liter.

The thresholds for the three detection procedures are adjusted so that the in-control ARL_{0s} are 10 days (960 samples). For the two procedures based on S^3T and the quadratic score statistic, the length of the sliding window is chosen as 12.5 hours (50 samples). Table 9 reports the average and standard error of detection delays obtained from 100 simulated spills. For spills with high intensity, all three methods achieve similar performance regarding

Table 9: Simulated expected detection delay in hours (numbers in parentheses are standard errors).

Spill Intensity	S^3T	Quadratic Score Statistic	T^2
low	38.285 (3.655)	45.822 (4.675)	52.959 (5.035)
medium	26.301 (1.679)	28.522 (1.873)	30.753 (2.192)
high	25.519 (1.697)	25.489 (1.667)	25.563 (1.860)

detection delay, as strong signals are easier to be detected. However, when the signal is relatively weak (low and medium spill intensity), the proposed detection statistic S^3T significantly outperforms the other two methods.

3.5 Conclusions

In this chapter, we propose a novel efficient score statistic S^3T to detect the emergence of a spatial-temporal signal from a noisy background in both the offline and online settings. The statistic captures the spatial and temporal correlation simultaneously and enjoys a relatively low computational cost. An accurate approximation for its false alarm rate is presented. Numerical results based on simulated data, real solar flare data, and a case study of water quality monitoring show that the proposed S^3T statistic has a clear advantage than existing methods.

CHAPTER IV

COMBINING CONSTRAINED BAYESIAN OPTIMIZATION AND SPATIO-TEMPORAL CHANGE-POINT DETECTION FOR SENSOR NETWORK DESIGN

The advances of sensor technology have enabled online monitoring for complicated systems. A sensor network consists of a group of sensors dispersed in space and collects multiple data streams in real time. Design of a sensor network answers not only where the sensors should be placed in the space but also how the data streams collected by the sensors should be processed and analyzed. In this chapter, we focus on designing a sensor network for water quality monitoring on a river system. The goal is to generate timely information regarding water quality and enable quick detections on undesired contamination events.

In practice, sensors are usually subject to random measurement error. [26] point out that decision making based on a sensor network is unreliable if the inaccuracy of data is not handled by a scientific approach. In water quality monitoring, the detection of contamination based on sensor data is a statistical problem. [26] adopt statistical process control (SPC) methods including a Shewhart chart and a CUSUM chart in sensor network design when measurement error exists. However, these approaches do not capture the spatial and temporal correlation in the sensor data. In Chapter 3, we propose the score statistic S^3T , which detects the emergence of a signal from noisy background. The S^3T statistic captures both spatial and temporal correlation in the data and hence is particularly good at detecting weak signals.

In this chapter, (i) we formulate the problem of sensor network design for water quality monitoring on river systems as a joint problem of constrained black-box function optimization and online statistical change-point detection; (ii) we propose a new algorithm called Confidence-Set based Constrained Bayesian Optimization (CSCBO), which provides a flexible framework to handle noisy black-box function constraints and is easy to implement;

(iii) we extend the algorithm to tackle with a challenge that arises specifically in the sensor network design problem: we adopt the Wasserstein similarity metric to enable CSCBO to solve problems with high-dimensional binary search space; and (iv) finally we combine CSCBO with the S^3T statistic to find robust sensor networks for the Altamaha river and show advantages of the proposed methods.

The rest of the chapter is organized as follows. Section 4.1 provides backgrounds on the constrained black-box function optimization problem and Bayesian optimization approach. Section 4.2 presents the proposed CSCBO algorithm. Section 4.3 presents the formulation and methodologies for sensor network design on river systems. Section 4.4 contains experimental results. Section 4.5 concludes the chapter.

4.1 Background

In this section, we provide general settings of the constrained black-box function optimization problems and a brief overview of the BO method.

4.1.1 Constrained black-box function optimization

Suppose a stochastic system of interest have $J+1$ performance metrics. One is considered as the primary performance metric and the rest are used as guardrail performance metrics. We seek to find the optimal input parameter of the system that leads to the best primary performance metric subject to nonnegative constraints on each guardrail performance metric. In mathematical notation, we consider the following constrained optimization problem,

$$\begin{aligned} \max_{\mathbf{x} \in \mathfrak{X}} \quad & f_0(\mathbf{x}) \\ \text{s.t.} \quad & f_j(\mathbf{x}) \geq 0, \quad \forall j \in \{1, \dots, J\} \end{aligned} \tag{41}$$

where \mathbf{x} is a D -dimensional vector of input variables, $\mathfrak{X} \subset \mathbb{R}^D$ is a bounded search space, $f_0 : \mathfrak{X} \rightarrow \mathbb{R}$ is the primary performance function and $f_j : \mathfrak{X} \rightarrow \mathbb{R}$, ($j \geq 1$) denote the j th guardrail performance function. Note that the search space \mathfrak{X} is assumed to have been discretized if the original domain of the functions is continuous, and hence, \mathfrak{X} is a discrete set and contains a finite number of points.

The functions f_j , $j = 0, \dots, J$ are considered to be unknown black-box functions that are expensive to evaluate and subject to noise. In other words, for any point $\mathbf{x} \in \mathfrak{X}$, $f_j(\mathbf{x})$

cannot be evaluated analytically. Usually $f_j(\mathbf{x})$, $j = 0, \dots, J$ are expectations of random estimates based on simulations for the performance metrics of the system given a parameter \mathbf{x} ,

$$f_j(\mathbf{x}) = \mathbb{E}[y_j(\mathbf{x})], j \in \{0, \dots, J\}.$$

Here $y_j(\mathbf{x})$ is a simulation estimate of the j th performance metric, which can be a single observation or a sample average of multiple observations. We assume that

$$y_j(\mathbf{x}) \sim \mathbb{N}\left(f_j(\mathbf{x}), \lambda_j(\mathbf{x})\right),$$

where $\lambda_j(\cdot)$ is referred to as the sampling variance of the j th performance metric. In practice, $\lambda_j(\cdot)$ is not necessarily known and should be estimated if unknown. From one run of simulation with input parameter \mathbf{x} , we are able to evaluate all performance metrics jointly and obtain a vector of observations $[y_0(\mathbf{x}), y_1(\mathbf{x}), \dots, y_J(\mathbf{x})]$.

4.1.2 Bayesian Optimization

Bayesian optimization (BO) is originally proposed to solve the unconstrained version of problem (41), i.e., $J = 0$. A BO algorithm is a sequential procedure and typically consists of two components: a surrogate model characterizing the function and an acquisition function guiding evaluations. In practice, Gaussian Process (GP) is most widely adopted to model the black-box objective function due to its flexibility and tractability. In the initial stage of the optimization, a GP prior is put over the function which is specified by a mean function $\mu_0(\mathbf{x}) : \mathfrak{X} \rightarrow \mathbb{R}$ and a kernel function $K_0(\mathbf{x}, \mathbf{x}') : \mathfrak{X} \times \mathfrak{X} \rightarrow \mathbb{R}$,

$$f_0(\mathbf{x}) \sim GP\left(\mu_0(\mathbf{x}), K_0(\mathbf{x}, \mathbf{x}')\right).$$

Suppose that we have evaluated f_0 at n points $\mathbf{x}^{(1:n)} := \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}\}$, and obtained corresponding observations $y_0^{(1:n)} := \{y_0^{(1)}, y_0^{(2)}, \dots, y_0^{(n)}\}$. Then, we can calculate the posterior distribution on f_0 by combining the prior and the observations based on the Bayes rule. The posterior is still a GP

$$f_0(\mathbf{x}) \Big| \mathbf{x}^{(1:n)}, y_0^{(1:n)} \sim GP\left(\mu_0^{(n)}(\mathbf{x}), K_0^{(n)}(\mathbf{x}, \mathbf{x}')\right),$$

where

$$\mu_0^{(n)}(\mathbf{x}) = \mu_0(\mathbf{x}) + K_0(\mathbf{x}, \mathbf{x}^{(1:n)}) \left(K_0(\mathbf{x}^{(1:n)}, \mathbf{x}^{(1:n)}) + \text{diag}\{\lambda(\mathbf{x}^{(1)}), \dots, \lambda(\mathbf{x}^{(n)})\} \right)^{-1} (y_0^{(1:n)} - \mu_0(\mathbf{x}^{(1:n)})), \quad (42)$$

and

$$K_0^{(n)}(\mathbf{x}, \mathbf{x}') = K_0(\mathbf{x}, \mathbf{x}') - K_0(\mathbf{x}, \mathbf{x}^{(1:n)}) \left(K_0(\mathbf{x}^{(1:n)}, \mathbf{x}^{(1:n)}) + \text{diag}\{\lambda(\mathbf{x}^{(1)}), \dots, \lambda(\mathbf{x}^{(n)})\} \right)^{-1} K_0(\mathbf{x}^{(1:n)}, \mathbf{x}'). \quad (43)$$

Typically, a BO algorithm determines the next evaluation point by maximizing an acquisition function which relies on the posterior distribution of f_0 ,

$$\mathbf{x}^{(n+1)} = \arg \max_{\mathbf{x} \in \mathfrak{X}} \text{acquisition_function}(\mathbf{x}).$$

Expected improvement (EI) [23] is one of the most commonly used acquisition function, which is calculated as follows,

$$\begin{aligned} \text{EI}(\mathbf{x}|f_0^*) &= \mathbb{E}[\max(0, f_0(\mathbf{x}) - f_0^*)] \\ &= (\mu_0^{(n)}(\mathbf{x}) - f_0^*)\Phi(Z) + \sqrt{K_0^{(n)}(\mathbf{x}, \mathbf{x})}\phi(Z), \end{aligned} \quad (44)$$

where $Z = \frac{\mu_0^{(n)}(\mathbf{x}) - f_0^*}{\sqrt{K_0^{(n)}(\mathbf{x}, \mathbf{x})}}$, f_0^* is the current best function value found and the expectation is taken over the posterior distribution of f_0 . Note that the original EI proposed in [23] assumes that $\lambda(\mathbf{x}) = 0$, $\forall \mathbf{x} \in \mathfrak{X}$, i.e., the function to be optimized is not subject to evaluation noise. When noise exists, calculating (44) is difficult as f_0^* is unknown. [2] propose to modify EI by replacing f_0^* in (44) with the GP posterior mean estimate of the best function value $\mu_0^{(n)*} = \max_{\mathbf{x} \in \mathfrak{X}} \mu_0^{(n)}(\mathbf{x})$,

$$\text{EI}(\mathbf{x}|\mu_0^{(n)*}) = (\mu_0^{(n)}(\mathbf{x}) - \mu_0^{(n)*})\Phi(Z) + \sqrt{K_0^{(n)}(\mathbf{x}, \mathbf{x})}\phi(Z). \quad (45)$$

For the rest of the chapter, we use EI to refer to the modified version defined in (45).

Another example of acquisition function is the upper confidence bound (UCB) [62]. UCB returns the next evaluation point with the highest upper confidence interval based on the posterior distribution of f_0 ,

$$\text{UCB}(\mathbf{x}, \alpha) = \mu_0^{(n)}(\mathbf{x}) + z_\alpha \sqrt{K_0^{(n)}(\mathbf{x}, \mathbf{x})}, \quad (46)$$

where z_α is the α th quantile of the standard normal distribution.

While EI and UCB are powerful methods and are successful in many of applications, they deal with unconstrained problems only. In the next section, we propose an efficient and practical algorithm that is capable of handling optimization problem with noisy black-box function constraints.

4.2 *Confidence-Set based Constrained Bayesian Optimization*

In this section, we present a new algorithm for problem (41) called Confidence-Set based Constrained Bayesian Optimization (CSCBO).

To solve the problem defined in (41), CSCBO puts independent GP priors over each function f_j at the initialization step,

$$f_j(\mathbf{x}) \sim GP\left(\mu_j(\mathbf{x}), K_j(\mathbf{x}, \mathbf{x}')\right), \quad j = 0, \dots, J, \quad (47)$$

where $\mu_j()$ and $K_j()$ are the prior mean and prior kernel functions for f_j . In the n th iteration, suppose the algorithm chooses to evaluate the functions at $\mathbf{x}^{(n)}$ and obtains a vector of observations $[y_0^{(n)}, y_1^{(n)}, \dots, y_J^{(n)}]$. The algorithm first updates the posterior distribution for each function independently,

$$f_j(\mathbf{x}) \sim GP\left(\mu_j^{(n)}(\mathbf{x}), K_j^{(n)}(\mathbf{x}, \mathbf{x}')\right), \quad j = 0, \dots, J, \quad (48)$$

where $\mu_j^{(n)}()$ and $K_j^{(n)}()$ denote the posterior mean and kernel functions for f_j after n evaluations and can be calculated using (42) and (43).

Next, we construct a *confidence set* (denoted as CS_1) using the posterior distributions, which is a subset of the search space \mathfrak{X} and eliminates points based on the criteria described in the following. We first define a reward function for each constraint as follows,

$$r_j(\mathbf{x}) = \begin{cases} 1, & \text{if } f_j(\mathbf{x}) \geq 0; \\ 0, & \text{otherwise,} \end{cases}$$

where $j = 1, \dots, J$. We then calculate the expected reward function with respect to the posterior distributions of the constraint functions,

$$\mathbb{E}[r_j(\mathbf{x})] = \mathbb{E}\left[\mathbb{1}\left\{f_j(\mathbf{x}) \geq 0\right\}\right] = \mathbb{P}\left(f_j(\mathbf{x}) \geq 0\right) = \Phi\left(\frac{\mu_j^{(n)}(\mathbf{x})}{\sqrt{K_j^{(n)}(\mathbf{x}, \mathbf{x})}}\right),$$

Algorithm 1 CSCBO

1: **Input:** GP priors $\mu_j()$, $K_j()$, $j = 0, \dots, J$; h_1 ; h_2 .
2: **Initialization:** set $n = 0$, $CS_1 = \mathfrak{X}$, and randomly pick $\mathbf{x}^{(0)} \in CS_1$.
3: **while** stopping criteria is not met **do**
4: evaluate at $\mathbf{x}^{(n)}$ and obtain $[y_0^{(n)}, y_1^{(n)}, \dots, y_J^{(n)}]$.
5: update GP posteriors $\mu_j^{(n)}()$ and $K_j^{(n)}()$ for $j = 0, \dots, J$.
6: set $CS_1 = \left\{ \mathbf{x} \in \mathfrak{X} : \Phi\left(\frac{\mu_j^{(n)}(\mathbf{x})}{\sqrt{K_j^{(n)}(\mathbf{x}, \mathbf{x})}}\right) \geq h_1 \ \forall j \in \{1, \dots, J\} \right\}$.
7: obtain $\mu_0^{(n)*} = \max_{\mathbf{x} \in CS_1} \mu_0^{(n)}(\mathbf{x})$ and $\mathbf{x}^{(n+1)} = \arg \max_{\mathbf{x} \in CS_1} \text{EI}(\mathbf{x} | \mu_0^{(n)*})$.
8: set $CS_2 = \left\{ \mathbf{x} \in \mathfrak{X} : \Phi\left(\frac{\mu_j^{(n)}(\mathbf{x})}{\sqrt{K_j^{(n)}(\mathbf{x}, \mathbf{x})}}\right) \geq h_2 \ \forall j \in \{1, \dots, J\} \right\}$.
9: report current optimal feasible solution $\mathbf{x}_*^{(n)} = \arg \max_{\mathbf{x} \in CS_2} \mu_0^{(n)}(\mathbf{x})$.
10: $n = n + 1$.
11: **end while**
12: **return** optimal feasible solution $\mathbf{x}_*^{(n)}$.

where $\mathbb{1}\{\cdot\}$ denotes an indicator function and $\Phi(\cdot)$ is the cumulative density function (CDF) of a standard normal distribution. Note that $f_j(\mathbf{x})$, $j = 1, \dots, J$ are treated as Gaussian random variables with mean $\mu_j^{(n)}(\mathbf{x})$ and variance $K_j^{(n)}(\mathbf{x}, \mathbf{x})$. CS_1 is constructed based on the expected rewards for all constraints,

$$CS_1 = \{\mathbf{x} \in \mathfrak{X} : \mathbb{E}[r_j(\mathbf{x})] \geq h_1, \ \forall j = 1, \dots, J\}, \quad (49)$$

where $h_1 \in [0, 1]$ is a pre-specified constant. Namely, CS_1 contains the points whose probabilities of satisfying each constraint are at least h_1 . In each iteration, the algorithm selects the next evaluation point from CS_1 based on the modified EI criterion defined in (45),

$$\mathbf{x}^{(n+1)} = \arg \max_{\mathbf{x} \in CS_1} \text{EI}(\mathbf{x} | \mu_0^{(n)*}). \quad (50)$$

For reporting optimal feasible solution found so far, another confidence set CS_2 is constructed with a different parameter h_2 . The point with the largest posterior mean in CS_2 is reported as the best feasible solution. The formal procedure of CSCBO is presented in Algorithm 1.

A tradeoff that one faces when dealing with a constrained black-box function optimization problem is whether to search for a promising solution in terms of the primary performance metric first and then check the feasibility or attempt to identify the boundary between feasible and infeasible regions first and then search for promising solutions. While

eventually an optimal feasible solution is desired, the two manners may lead to different search schemes. We design CSCBO based on the idea that improving the primary performance metric should be of higher priority. Specifically, the feasibility of a solution matters only if the solution is likely to have a large value of the primary performance metric. Thus, in CSCBO, the acquisition function itself is purely based on the primary performance function f_0 , and the constraints only affect the search procedure via the confidence set CS_1 . CS_1 restricts the region where the next evaluation point is chosen and is supposed to only eliminate the points with strong evidence to be infeasible. Thus, h_1 should be a relatively small value. CS_1 is initialized to include all points in \mathfrak{X} and hence the restriction on the search space is small in the early iterations of the optimization. In practice, computation budget is not unlimited and the optimization algorithm is subject to termination within finite number of iterations. The confidence set CS_2 is constructed for the purpose of reporting an optimal feasible solution after a finite number of iterations. The optimal feasible solution identified in the n th iteration, which is denoted as $\mathbf{x}_\star^{(n)}$, is the point with the highest posterior mean in CS_2 . To gain more confidence in the feasibility of $\mathbf{x}_\star^{(n)}$, we recommend that h_2 is set as a relatively large value, and hence there is strong evidence for points in CS_2 to be feasible. A simulation study for different choices of h_1 and h_2 is presented in Section 4.2.2.

Stopping criteria. One may choose to terminate the algorithm after a pre-specified maximum number of iterations. Alternatively, we can terminate the algorithm after the same point has been reported as the optimal feasible solution in k consecutive iterations. A recommended value is $k = 50$.

4.2.1 Connection to confidence bounds

An alternative way to construct a confidence set is based on confidence bounds of a GP. We define an α -level confidence bound of the posterior GP for function f_j after n evaluations as follows,

$$CB_j(\mathbf{x}, \alpha) = \mu_j^{(n)}(\mathbf{x}) + z_\alpha \sqrt{K_j^{(n)}(\mathbf{x}, \mathbf{x})},$$

where z_α is the α th upper quantile of a standard normal distribution, i.e., $\mathbb{P}(Z \geq z_\alpha) = \alpha$ and $Z \sim \mathbb{N}(0, 1)$. Then, we can construct CS_1 as follows,

$$CS_1 = \left\{ \mathbf{x} \in \mathfrak{X} : CB_j(\mathbf{x}, \alpha_1) \geq 0 \ \forall j \in \{1, \dots, J\} \right\}, \quad (51)$$

where $\alpha_1 \in [0, 1]$ is a constant. In the following, we demonstrate that the two confidence sets defined in (49) and (51) are equivalent if $\alpha_1 = h_1$. First we have,

$$\mathbb{P}\left(f_j(\mathbf{x}) \geq CB_j(\mathbf{x}, \alpha_1)\right) = \mathbb{P}\left(f_j(\mathbf{x}) \geq \mu_j^{(n)}(\mathbf{x}) + z_{\alpha_1} \sqrt{K_j^{(n)}(\mathbf{x}, \mathbf{x})}\right) = \alpha_1.$$

Also, we have that $\mathbb{P}(f_j(\mathbf{x}) \geq t)$ is a monotone decreasing function of t as $f_j(\mathbf{x})$ is treated as a random variable. Therefore,

$$CB_j(\mathbf{x}, \alpha_1) \geq 0 \Leftrightarrow \mathbb{E}[r_j(\mathbf{x})] = \mathbb{P}\left(f_j(\mathbf{x}) \geq 0\right) \geq \alpha_1 = h_1.$$

4.2.2 Choices for h_1 and h_2

Here we demonstrate the impact of h_1 and h_2 values on the performance of CSCBO. We solve the following problem,

$$\begin{aligned} \max_{\mathbf{x} \in [0, 100]} \quad & \mathbb{E} \left[\frac{\sin^6(0.05\pi x)}{2^{2(\frac{x-10}{80})^2}} + \epsilon_0 \right] \\ \text{s.t.} \quad & \mathbb{E} \left[\sin(0.1\pi x) + H + \epsilon_1 \right] \geq 0, \end{aligned}$$

where H is a constant and ϵ_0 and ϵ_1 are independent normal random variables with mean 0 and variance σ_ϵ^2 . We set $H \in [-0.8, -0.99]$ to create scenarios where the proportion of feasible solution is relatively small ($H = -0.8$) and extremely small ($H = -0.99$). σ_ϵ is set as 0 or 0.01 to test the performance of CSCBO on deterministic and stochastic functions, respectively. The search space is uniformly discretized into 1000 points. We define a utility function $u(\mathbf{x})$ which is equal to $f_0(\mathbf{x})$ if \mathbf{x} is feasible and otherwise the worst objective function value achievable in the search space. In each iteration, we record the utility of the current optimal feasible solution reported by the algorithm $u(\mathbf{x}_\star^{(n)})$ and compute the utility gap $|u(\mathbf{x}_\star^{(n)}) - u(\mathbf{x}_\star)|$, where \mathbf{x}_\star is the true optimal feasible solution of the problem. The algorithm is initialized with a random point in the search space.

Figure 15 shows the mean of the utility gap for 500 replications with different randomized initialization. We can observe that the algorithm performs better if h_1 is set to a

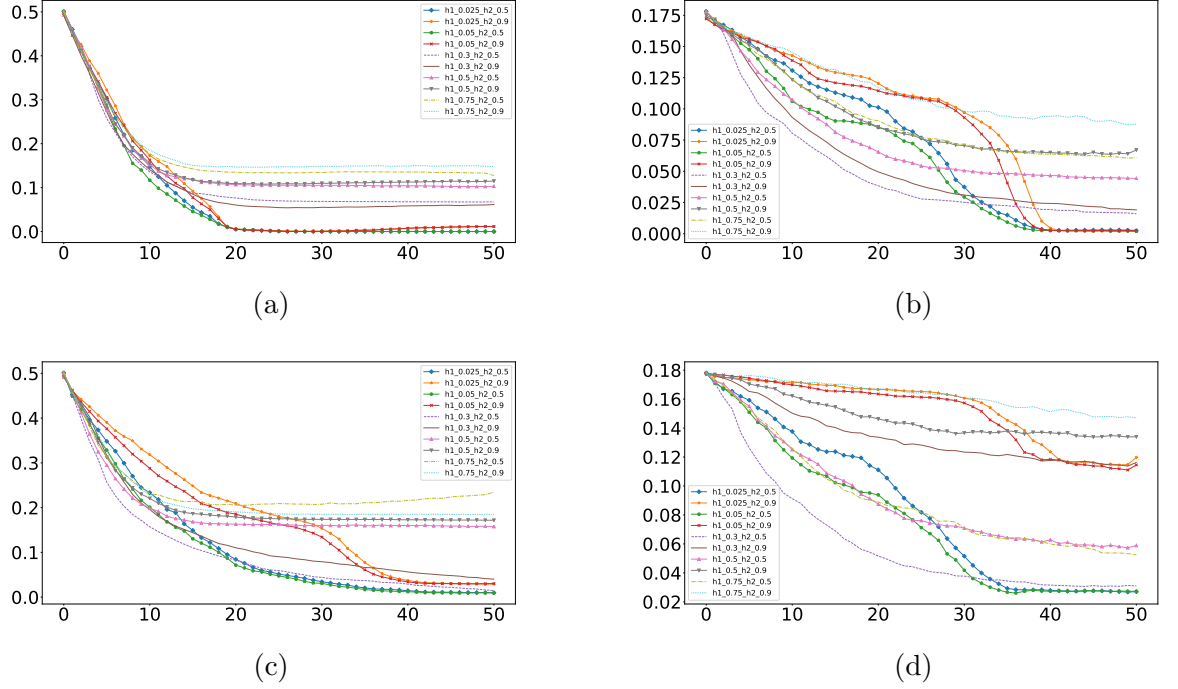


Figure 15: Mean utility gap between the true optimal feasible solution and the solution found by CSCBO: (a) deterministic, $H = -0.8$, (b) deterministic, $H = -0.99$, (c) stochastic, $H = -0.8$ and (d) stochastic, $H = -0.99$.

smaller value (e.g., $h_1 = 0.025$ or $h_1 = 0.05$), which achieves smaller mean utility gap when converging, than a larger value ($h_1 \in [0.3, 0.5, 0.75]$). The reason is that a large h_1 may restrict the search space too much in early iterations of the optimization when only a few function evaluations can be used to update the posterior distribution of the GPs, which leads to a lack of exploration. Therefore, h_1 should be set to a small value to avoid this issue. We also find that the algorithm converges quicker with $h_2 = 0.5$ than $h_2 = 0.9$. Such an advantage is more obvious for problems with lower proportion of feasible solutions ($H = -0.99$). This is due to that $h_2 = 0.9$ is sometimes too strict for the algorithm to find any feasible solution when the true proportion of feasible solutions is small. Based on these observations, we should avoid using a large h_2 value in practice.

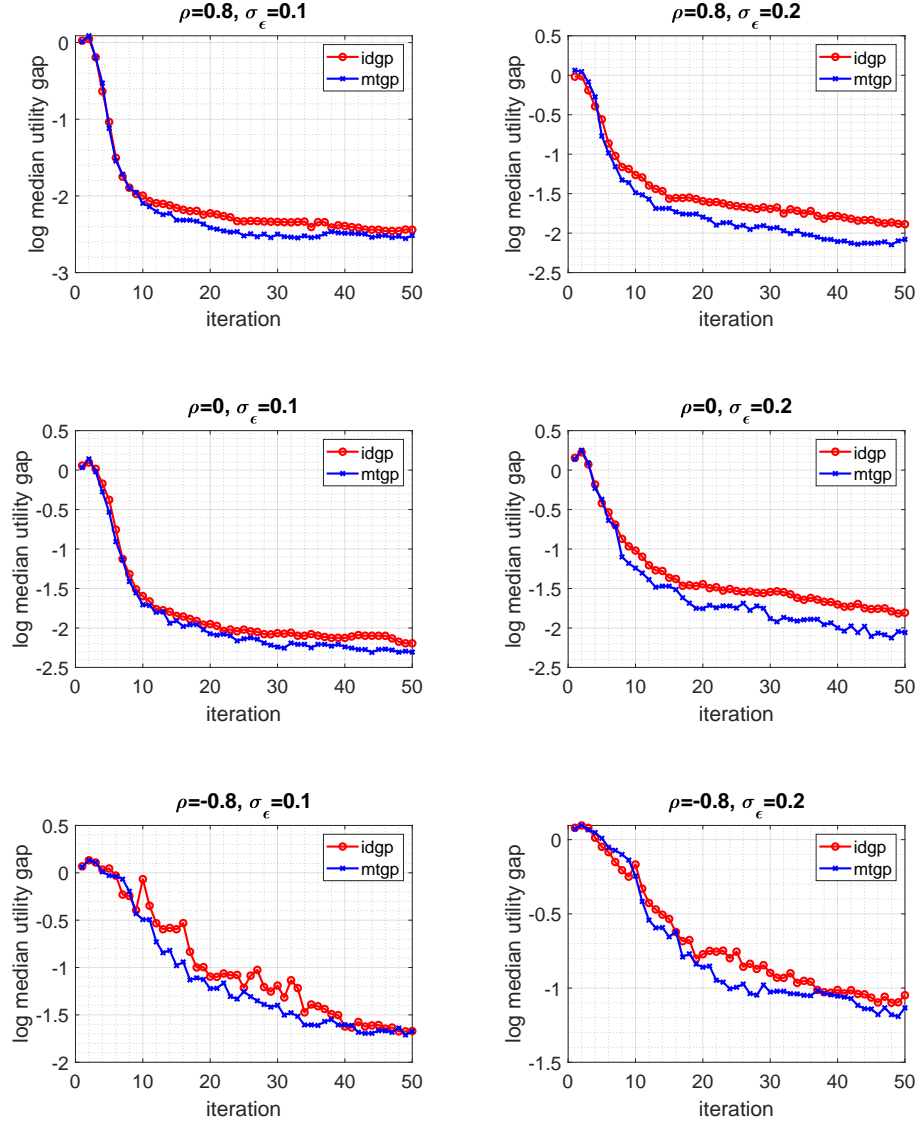


Figure 16: Comparison on log median utility gap between CSCBO with independent GPs (idgp) and multi-task GPs (mtgdp).

4.2.3 Extend to multi-task Gaussian processes (MTGP)

[4] extends the Gaussian processes to the case of multi-output functions by defining a covariance function $K_{\text{multi}}((\mathbf{x}, t), (\mathbf{x}', t'))$ between functions,

$$K_{\text{multi}}((\mathbf{x}, t), (\mathbf{x}', t')) = K_t(t, t') \otimes K_x(\mathbf{x}, \mathbf{x}'),$$

where K_x is the covariance function between input points and K_t is the covariance between functions. The posterior distribution of a multi-task GP can be obtained by the common approach once the covariance function K_{multi} is given. The benefit of MTGP is that information about multiple functions can be shared in the joint model. To take advantage of such benefit, we extend CSCBO by modeling all functions $f_j, j = 0, \dots, J$ using a MTGP and jointly obtain the posterior mean and variances for each function in every iteration. The rest steps proceed in the same way.

In the following, we demonstrate the performance of CSCBO with MTGP using synthetic problems in the form of (41) with $J = 1$. The search space is 1000 uniformly discretized points in $[0, 1]$. In each problem, an objective function and a constraint are randomly generated by a zero-mean MTGP with a squared exponential kernel of unit amplitude and length scale $\ell = 0.1$. The correlation between the two GPs used for generating objective functions and constraints is set to -0.8, 0 and 0.8. The sampling standard deviation σ_ϵ is set to 0.1 and 0.2. Note that in the experiment, the true information of the GP models used for generating functions is not given to the algorithm. We define the utility function $u(\mathbf{x})$ in the same way as in Section 4.2.2. For each setting, we randomly generate 50 problems and repeat the solving process 10 times with random initialization for each problem. Figure 16 presents the median utility gap on logarithmic scale for CSCBO based on MTGP and independent GPs. We can observe that CSCBO with MTGP outperforms independent GPs in all cases, although the difference is small when the sample variance is small. This indicates that MTGP indeed helps to improve the efficiency of CSCBO by capturing the correlation between the objective function and constraint. Interestingly, the advantage of MTGP exists even when $\rho = 0$. The reason is that although independently, the objective function and the constraint are randomly generated by GPs with same hyperparameters and hence tends to behave similarly.

4.3 Combine CSCBO and $\mathbf{S}^3\mathbf{T}$ for Sensor Network Design

In this section, we combine CSCBO and $\mathbf{S}^3\mathbf{T}$ to solve the problem of sensor network design for river water quality monitoring. The goal of the network is to quickly and accurately

detect contaminant events on the river.

4.3.1 Formulation

We first formulate the problem of optimal sensor placement as a constrained black-box function optimization problem as defined in (41) under the assumption that sensors are error-free, i.e., sensor measurement error does not exist. A river network is assumed to have D nodes, indexed from 1 to D and each one is a potential location to place a sensor. Denote the number of sensors as M and assume $M < D$. The decision variable \mathbf{x} is a set of nodes where the M sensors are placed. \mathbf{x} is a D -dimensional binary vector with $x_i = 1$ if a sensor is placed on the i th node and 0 otherwise. The search space $\mathfrak{X} = \{\mathbf{x} \in BV^D : |\mathbf{x}|_1 = M\}$, where BV^D denotes a D -dimensional binary vector space and $|\cdot|_1$ denotes the L_1 norm of a vector.

The sensors collect concentration data in real-time. At time t , the M sensors obtain a vector of concentration measurements $\mathbf{c}_t = [c_{1t}, \dots, c_{Mt}]$, where c_{it} denotes the data measured by the i th sensor at time t . Under the error-free assumption, the concentration measurements are collected with 100% accuracy. A detection statistic W_t is constructed based on the concentration measurements collected at and prior to t . A simple example of a detection statistic is

$$W_t = \max(c_{1t}, \dots, c_{Mt}). \quad (52)$$

The monitoring system will raise an alarm at time t warning a contaminant is detected if W_t exceeds an pre-specified threshold b . Since sensors are assumed to be error-free, if an alarm is raised, a true contaminant spill is detected successfully. Denote t_0 as the starting time of a contaminant spill event and $t_a(\mathbf{x})$ as the time stamp when the sensor network at \mathbf{x} raises an alarm,

$$t_a(\mathbf{x}) = \inf \left\{ t : W_t \geq b \right\}.$$

Then we define detection delay $T(\mathbf{x})$ as follows,

$$T(\mathbf{x}) = t_a(\mathbf{x}) - t_0. \quad (53)$$

$T(\mathbf{x})$ is the amount of time elapsed between the start and the detection of a contaminant event. Note that it is possible that a sensor network fails to detect the contaminant event.

In such a case, the sensor network will never raise an alarm and $T(\mathbf{x}) = \infty$. We define an indicator function

$$R(\mathbf{x}) = \begin{cases} 0, & \text{if the sensor network fails to detect a contaminant event, (i.e., } T(\mathbf{x}) = \infty); \\ 1, & \text{otherwise.} \end{cases} \quad (54)$$

Note that both $T(\mathbf{x})$ and $R(\mathbf{x})$ are random variables and can only be observed via stochastic simulations.

Adopting the same formulation in [41], we formulate the problem of optimal sensor placement for river water quality monitoring network with error-free sensors as follows,

$$\begin{aligned} \min_{\mathbf{x} \in \mathcal{X}} \quad & \mathbb{E}[T(\mathbf{x}) | R(\mathbf{x}) = 1] \\ \text{s.t.} \quad & \mathbb{E}[R(\mathbf{x})] \geq q, \end{aligned} \quad (55)$$

where $0 \leq q \leq 1$. We refer to $\mathbb{E}[T(\mathbf{x}) | R(\mathbf{x}) = 1]$ as the conditional expected detection delay and $\mathbb{E}[R(\mathbf{x})]$ as the reliability, both of which are stochastic black-box functions.

4.3.2 Measurement error

In practice, sensors are usually subject to measurement error, which means the concentration measurements collected by the sensors fluctuate around the true values by small random amounts. We assume that when there is no contaminant event, the concentration measurements collected by the M sensors at each time stamp are a series of i.i.d. multivariate normal random variables,

$$\mathbf{c}_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}), \quad t = 1, 2, \dots,$$

where $\mathbf{0}$ represent a M -dimensional 0 vector and $\mathbf{\Sigma}$ is the spatial covariance matrix of the measurements. We often have enough reference data in the in-control case (no contaminant event) to estimate the mean and temporal correlation and to whiten the data, hence assuming \mathbf{c}_t has zero mean and no temporal correlation when no contaminant event is reasonable. We also assume that $\mathbf{\Sigma}$ is known or has been estimated. When a contaminant event occurs, the mean of \mathbf{c}_t is equal to the true concentration level of the contaminant in the water. In addition, temporal correlation exists due to the contaminant transport along the river.

In this case, a sensor network needs to detect a contaminant event with not only a low detection delay and a high reliability but also a low false alarm rate. A false alarm is defined as the case when there is no contaminant event but the sensors raise an alarm. We use the in-control average-run-length as the metric to quantify the frequency of a false alarm, which is the average number of time steps until an alarm is raised where there is no contaminant. For a specific \mathbf{x} , we use $\text{ARL}_0(\mathbf{x})$ to denote the in-control average-run-length of the sensor network at \mathbf{x} . Then, the optimal sensor placement problem in the presence of measurement error is formulated as follows,

$$\begin{aligned} \min_{\mathbf{x} \in \mathfrak{X}} \quad & \mathbb{E}[T(\mathbf{x}) | R(\mathbf{x}) = 1] \\ \text{s.t.} \quad & \mathbb{E}[R(\mathbf{x})] \geq q, \quad q \in [0, 1], \\ & \text{ARL}_0(\mathbf{x}) \geq \text{ARL}_{\text{target}}. \end{aligned} \tag{56}$$

The detection delay $T(\mathbf{x})$ not only depends on the sensor placement \mathbf{x} but the detection statistic W_t as well. In this chapter, we use the Score Statistic for Spatio-Temporal surveillance (S^3T) proposed in Chapter 3 as the detection statistic. S^3T statistic captures both spatial and temporal correlation information and is good at detecting weak signals. At each time stamp t , a sliding window of length ω is constructed to form the statistic. Define the following notations,

$$\mathbf{c}_{(t-\omega+1:t)} = [\mathbf{c}_{t-\omega+1}^\top, \dots, \mathbf{c}_t^\top],$$

$$\mathbf{\Sigma}_\omega = \mathbf{I}_\omega \otimes \mathbf{\Sigma},$$

$$\mathbf{V}_\omega(\theta) = \mathbf{R}_\omega(\theta) \otimes \mathbf{\Sigma},$$

$$\mathbf{A}_\omega(\theta) = \mathbf{\Sigma}_\omega^{-1} \mathbf{V}_\omega(\theta),$$

$$\mathbf{B}_\omega(\theta) = \mathbf{\Sigma}_\omega^{-1/2} \mathbf{V}_\omega(\theta) \mathbf{\Sigma}_\omega^{-1/2},$$

where \mathbf{I}_ω is a ω -by- ω identity matrix, \otimes denotes the Kronecker product and the matrix $\mathbf{R}_\omega(\theta) \in \mathbb{R}^{\omega \times \omega}$ captures the temporal correlation of the concentration measurements collected by the sensors when a contaminant event occurs. As suggested in [8] and [9], we assume that the temporal correlation follows a first-order vector autoregressive VAR(1)

model and hence $[\mathbf{R}_\omega(\theta)]_{i,j} = \theta^{|i-j|}, \forall i, j \in \{1, \dots, \omega\}$. Then, the $\mathbf{S}^3\mathbf{T}$ statistic with window length ω is calculated as follows,

$$W_t(\omega) = \max_{\theta \in \Theta} \frac{\mathbf{c}_{(t-\omega+1:t)}^\top \boldsymbol{\Sigma}_\omega^{-1} \mathbf{V}_\omega(\theta) \boldsymbol{\Sigma}_\omega^{-1} \mathbf{c}_{(t-\omega+1:t)} - \text{tr}(\mathbf{A}_\omega(\theta))}{\sqrt{2\text{tr}(\boldsymbol{\Sigma}_\omega^{-1} \mathbf{V}_\omega(\theta) \boldsymbol{\Sigma}_\omega^{-1} \mathbf{V}_\omega(\theta))}}, \quad (57)$$

where $\text{tr}(\cdot)$ denotes the trace of a matrix and Θ is a pre-specified parameter set.

To achieve a desired ARL_0 target, the threshold b needs to be adjusted. When spatial and temporal correlations exist in the data and $\mathbf{S}^3\mathbf{T}$ is used as the detection statistic, b may take different values for different \mathbf{x} as the spatial and temporal covariances depend on the locations of the sensors. In this case, adjusting b is needed in every iteration of the optimization process. In practice, b can be adjusted based on simulations if the distribution of in-control data (i.e., when no contaminant event) is known or can be estimated. In the following, we briefly introduce an efficient simulation approach to adjusting b which is based on the fact that the distribution of run-length (denoted as RL, which is the number of time steps until a false alarm is raised) approximately follows exponential distribution. We first simulate m sequences of in-control data with a fixed length n , $\{\mathbf{c}_{i1}, \dots, \mathbf{c}_{in}\}_{i=1}^m$, where $\mathbf{c}_{ij} \in \mathbb{R}^M$, and construct m sequences of detection statistics based on the simulated data, $\{W_{i1}, \dots, W_{in}\}_{i=1}^m$. Next we obtain the maximum value of the detection statistics in each sequence, $\{W_i^*\}_{i=1}^m$, where $W_i^* = \max(W_{i1}, \dots, W_{in})$. For a specific b , we calculate $\hat{p} = \frac{\sum_{i=1}^m \mathbb{1}\{W_i^* > b\}}{m}$, where $\mathbb{1}\{\cdot\}$ denotes an indicator function. Then, we have the following approximation,

$$\hat{p} \approx \mathbb{P}(\text{RL} > n) \approx e^{-\lambda n},$$

where $\lambda = \mathbb{E}[\text{RL}] = \text{ARL}_0$. We can now obtain an approximated ARL_0 if a specific b value is used as the threshold. For a target ARL_0 , we can find the appropriate b using a binary search. Such an approach is efficient because n needs not to be very large. We can obtain an accurate b using $n = \frac{\text{ARL}_{\text{target}}}{10}$.

4.3.3 Process simulation

In order to evaluate $E[T(\mathbf{x})|R(\mathbf{x}) = 1]$ and $E[R(\mathbf{x})]$ for different sensor placements on a river network, a simulation model that is capable of simulating hydrodynamics and contaminant transport is needed. The Storm Water Management Model (SWMM, [49]) developed by the United States Environmental Protection Agency is widely used in environmental engineering for water-related study. SWMM requires geologic, geometric and fundamental hydrodynamics data to construct a river network.

The random contaminant events are simulated by the SWMM model. The randomness of a contaminant event is due to the randomness of the location, intensity, duration and starting time of the spill as well as random rain events. For a specific sensor placement \mathbf{x} , a SWMM run is implemented as follows:

Step 1: Randomly generate the location, intensity, duration and starting time of a contaminant spill and rain pattern.

Step 2: Simulate contaminant transport along the river in consideration with the input data generated from Step 1.

Step 3: Evaluate $T(\mathbf{x})$ and $R(\mathbf{x})$ based on the concentration measurements collected at \mathbf{x} .

4.3.4 Wasserstein similarity metric

Bayesian optimization algorithms take advantage of the “smoothness” of a function, i.e., the function has similar values for similar decision variables. The similarities among different points in the search space are captured by the kernel function of the GP. An example of the commonly used kernel functions is the square-exponential kernel,

$$K(\mathbf{x}, \mathbf{x}') = \phi^2 \exp \left\{ -\frac{1}{2} d(\mathbf{x}, \mathbf{x}'; \ell)^2 \right\},$$

where $\phi^2 > 0$ is the variance parameter and $d(\cdot, \cdot; \ell) \geq 0$ is a similarity metric parametrized by the length scale parameter $\ell > 0$. Smaller $d(\mathbf{x}, \mathbf{x}'; \ell)$ indicates higher similarity between \mathbf{x} and \mathbf{x}' in performance. If the decision variables are continuous variables or ordinal discrete

variables, a natural choice for d is the Euclidean distance,

$$d(\mathbf{x}, \mathbf{x}'; \ell)^2 = \sum_{i=1}^D \frac{(\mathbf{x}_i - \mathbf{x}'_i)^2}{\ell}. \quad (58)$$

However, the decision variable \mathbf{x} is a D -dimensional binary vector in problems (55) and (56), which makes Euclidean distance an inappropriate choice for the similarity metric. The reason is demonstrated in the following example. Figure 17 shows three different sensor placements (number of sensors $M = 2$) on a hypothetical river with $D = 6$ nodes. Denote the three sensor placements in Figure 17 (a), (b) and (c) as \mathbf{x}_a , \mathbf{x}_b and \mathbf{x}_c , respectively. We have $\mathbf{x}_a = [1, 0, 0, 1, 0, 0]$, $\mathbf{x}_b = [0, 1, 1, 0, 0, 0]$ and $\mathbf{x}_c = [0, 0, 0, 0, 1, 1]$. The Euclidean distance between \mathbf{x}_a and \mathbf{x}_b is 2 which is equal to the Euclidean distance between \mathbf{x}_a and \mathbf{x}_c . However, as we can observe from Figure 17, from \mathbf{x}_a to \mathbf{x}_b , we only need to move the two sensor from Node 1 to 2 and from Node 4 to 3 by small amounts, which will not greatly affect the performance. Hence, \mathbf{x}_b is a much more similar placement to \mathbf{x}_a than \mathbf{x}_c . Such difference is not captured by the Euclidean distance metric.

We propose to use the Wasserstein metric as the similarity metric. It is originally proposed as a measure of discrepancy between two distributions and is widely applied in image processing [50]. We adjust the metric to the setting of sensor placement on a river network. For a specific sensor placement, we first convert the D -dimensional binary vector \mathbf{x} into a series of weights $\{w_i > 0\}_{i=1}^D$ by assigning $w_i = \frac{1}{M}$ if a sensor is placed at the i th node and $w_i = 0$ otherwise. Note that $\sum_{i=1}^D w_i = 1$. For two different sensor placements \mathbf{x} and \mathbf{x}' , and the corresponding weights $\{w_i > 0\}_{i=1}^D$ and $\{v_i > 0\}_{i=1}^D$, the Wasserstein distance, denoted as $d_w(\mathbf{x}, \mathbf{x}')$, is intuitively the minimum amount of work required to move all sensors from placement \mathbf{x} to \mathbf{x}' . Formally, $d_w(\mathbf{x}, \mathbf{x}')$ is equal to the optimum of the following linear program [50]:

$$\begin{aligned} \min \quad & \sum_{i=1}^M \sum_{j=1}^M g_{ij} C_{ij} \\ \text{s.t.} \quad & g_{ij} \geq 0, \quad 1 \leq i \leq M, 1 \leq j \leq M, \\ & \sum_{j=1}^M g_{ij} = w_i, \quad 1 \leq i \leq M, \\ & \sum_{i=1}^M g_{ij} = v_j, \quad 1 \leq j \leq M, \\ & \sum_{i=1}^M \sum_{j=1}^M g_{ij} = 1, \end{aligned} \quad (59)$$

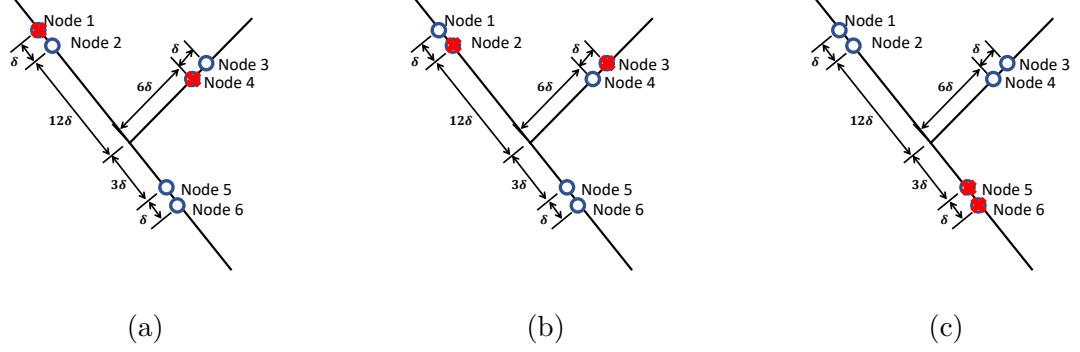


Figure 17: Three different sensor placements (number of sensors $M = 2$) on a hypothetical river with $D = 6$ nodes. Sensors are marked by the red crosses. The stream distances between each node are also marked on the plots.

where g_{ij} can be intuitively understood as the number of sensor moved from node i to j and C_{ij} is the cost of moving one sensor from node i to j . We use the stream distance, namely the shortest distance along the river between node i and j as C_{ij} . For the previous example, $d_w(\mathbf{x}_a, \mathbf{x}_b) = 2\delta$ and it is much smaller than $d_w(\mathbf{x}_a, \mathbf{x}_c) = 26\delta$.

4.4 Experiments

In this section, we solve problems (55) and (56) for the Altamaha River using the proposed methods. Backgrounds on the Altamaha River and SWMM model can be found in Section 2.5.1.

4.4.1 Simulation setup

In each SWMM run, 100 contaminant events are generated and each event is a single instantaneous spill at a different node. Note that the spills generated by the same SWMM run share the same rain pattern, but have different intensity and starting time. For a specific sensor placement \mathbf{x} , the output of each SWMM run are estimates of the conditional expected detection delay and reliability based on the 100 contaminant events.

4.4.2 Error-free sensors

We first assume that the sensors are error-free and consider problem (55).

4.4.2.1 Settings

In this case, false alarms are not of our concern and the threshold b is chosen arbitrarily. We adopt the same setting in [41]: $b \in \{0.05, 0.0001\}$ (milligram/litter). Input parameters for CSCBO include h_1 and h_2 which control the way CS_1 and CS_2 are constructed, as well as k in the stopping criteria which is the minimum number of consecutive iterations where the same point is reported as the optimal feasible solution. We set $h_1 = 0.025$, $h_2 = 0.5$ and $k = 50$. We use independent GPs in the algorithm, as the benefit of MTGP is not obvious due to the high-dimensional nature of the problem. A square-exponential kernel as defined in (58) is used as the kernel function of the GP for $E[R(\mathbf{x})]$. The kernel function of the GP for $E[T(\mathbf{x})|R(\mathbf{x}) = 1]$ is a Matérn 5/2 model, as defined in the following,

$$K(\mathbf{x}, \mathbf{x}') = \phi^2 \left(1 + \frac{\sqrt{5}d_w(\mathbf{x}, \mathbf{x}')}{\ell} + \frac{5d_w(\mathbf{x}, \mathbf{x}')^2}{3\ell^2} \right) \exp \left\{ -\frac{\sqrt{5}d_w(\mathbf{x}, \mathbf{x}')}{\ell} \right\},$$

where $\ell > 0$ is the length scale parameter. Note that we do not need to specify the prior mean, prior variance and length scale parameter ℓ for the GPs as the maximum likelihood estimates of these parameters [45] are used. The minimum reliability level is set as $q = 0.9$. We test two possible values for the number of sensors, $M \in \{5, 7\}$.

Implementing a Bayesian optimization algorithm is challenging when the size of the search space \mathfrak{X} is large. In our problem, the size of \mathfrak{X} grows exponentially in the number of possible sensor location on the river network D when the number of sensors M is fixed. For example, if $M = 5$ and the number of possible sensor location $D = 100$, we have $|\mathfrak{X}| = \binom{100}{5} = 75287520$, which makes implementing CSCBO on a personal desktop impossible. To speed up the computation, we exclude the nodes on the most upstream points in the search as placing sensors on these nodes is obviously not beneficial. We limit the number of possible sensor locations to $D = 50$ for $M = 5$ and $D = 30$ for $M = 7$.

In the experiment, we compare CSCBO with NP + PFM [41]. Three performance measures are used in the comparison: (i) estimated conditional expected detection delay (ECEDD) and (ii) estimated reliability (ER) of the optimal feasible solution found by the algorithm as well as (iii) the number of SWMM runs required until termination of the algorithm (NUM).

Table 10: Performance metrics of the optimal feasible solutions found by CSCBO and NP + PFM. The unit of ECEDD is hour.

M	b	CSCBO			NP + PFM		
		ECEDD	ER	NUM	ECEDD	ER	NUM
5	0.05	59.16	0.901	205	55.54	0.902	614
5	0.0001	49.01	0.930	95	46.13	0.930	154
7	0.05	45.49	0.903	237	44.47	0.916	800
7	0.0001	38.48	0.930	113	37.05	0.930	205

4.4.2.2 Results

Figure 18 shows the optimal feasible solution found by CSCBO in circles and NP + PFM in triangles for different settings. Table 10 presents the corresponding performance metrics of the two methods. From the results, we observe that all the solutions found by CSCBO are feasible, indicating confidence-set is a effective method to deal with optimization problems with black-box function constraints. In addition, CSCBO is able to find competitive sub-optimal solutions using much fewer number of SWMM runs than NP + PFM. Considering that NP + PFM visits multiple solutions in each iteration (200 solutions per iteration in the experiment), CSCBO converges much more quickly than NP + PFM in terms of the number of solutions visited. Such advantage of CSCBO leads to tremendous savings on computation resources considering SWMM is a highly computational-intensive model and each SWMM run takes around 2 hours on a personal desktop with Intel Core I5 CPU.

It turns out that CSCBO tends to converge to a local optimum in problem (55). We find that the algorithm has difficulty finding solutions that are better than the ones reported in Figure 18 when the algorithm is allowed to run more iterations. The difficulty mainly arises from the high dimensional nature of the problem. BO algorithms usually perform well on problems with low dimensionality and continuous decision variables but suffer in high-dimensional problems with nominal variables. Although the Wasserstein similarity metric discussed in Section 4.3.4 enables CSCBO to deal with a problem with D -dimensional binary variables, it is still difficult for the algorithm to identify the true optimal feasible solution given the huge search space. Hence, CSCBO is a favorable method when computation

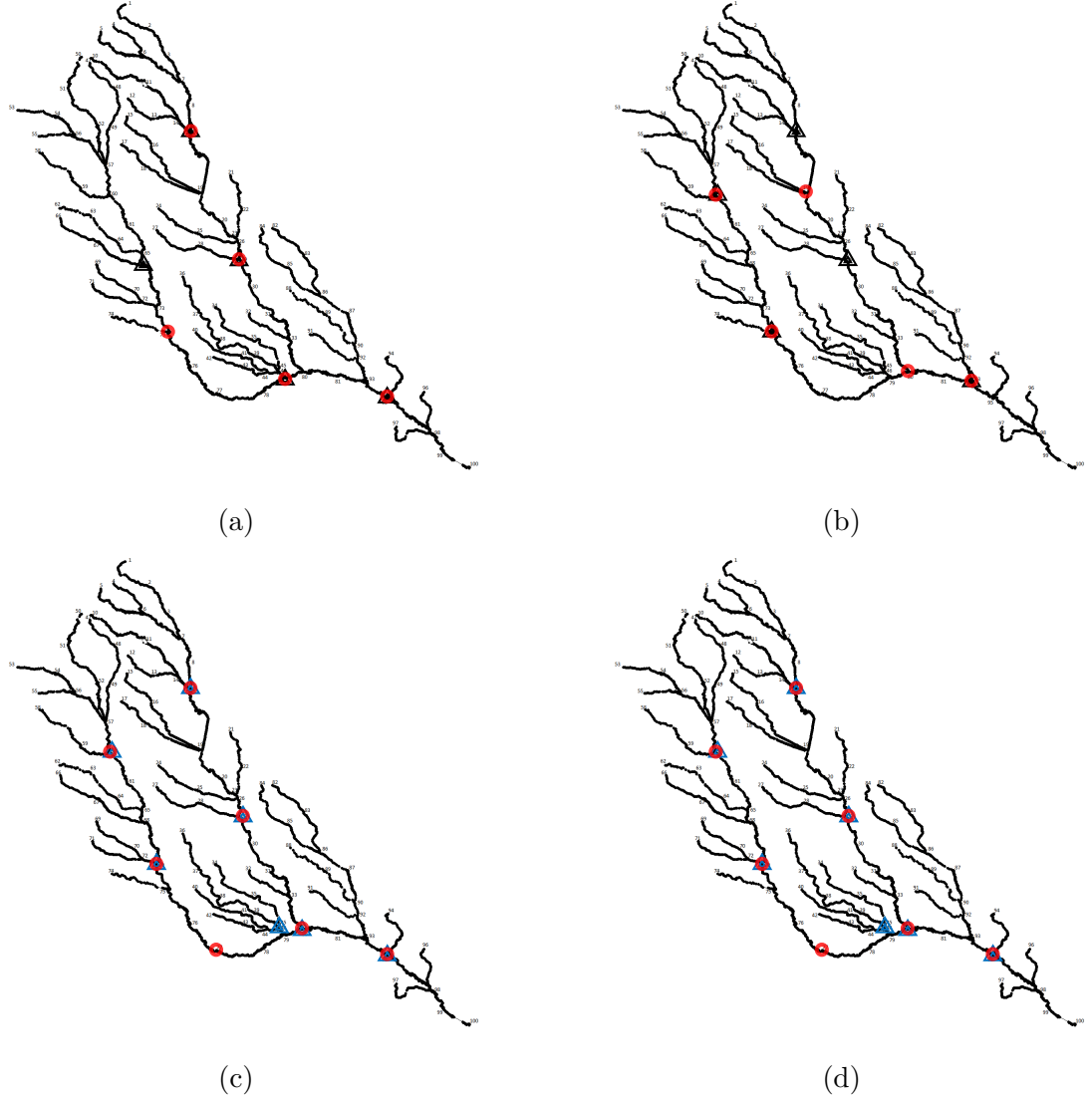


Figure 18: Optimal feasible solutions found by CSCBO (circle) and NP + PFM (triangle): (a) $M = 5$ and $b = 0.05$, (b) $M = 5$ and $b = 0.0001$, (c) $M = 7$ and $b = 0.05$ and (d) $M = 7$ and $b = 0.0001$.

budget is tight and a local optimal feasible solution with little sacrifice on the objective function is acceptable.

4.4.3 Sensors with measurement error

Here we assume that sensor measurement error exists and apply CSCBO and S^3T statistic to solving problem (56).

Table 11: ECEDD (in hours) of sensor placements marked by circles in Figure 18(a) and (c) using different detection statistics. Numbers in parentheses are standard errors.

M	ζ_1	S^3T	Shewhart	max-CUSUM
5	0.005	71.13 (0.21)	81.70 (0.34)	73.21 (0.20)
5	0.01	72.91 (0.24)	85.18 (0.39)	75.65 (0.23)
7	0.005	56.48 (0.20)	66.77 (0.31)	58.43 (0.16)
7	0.01	58.14 (0.22)	75.59 (0.41)	61.57 (0.22)

4.4.3.1 Settings

We adopt the “tail-up” model proposed in [69] (see 2.5.2 for details) with an exponential correlation function for the spatial correlation of the concentration measurements collected by the sensors. The concentration measurements follow the same spatial correlation structure before and after a contaminant event. The maximum likelihood estimate $\hat{\zeta}_2 = 0.68$ is obtained based on the data simulated by the SWMM model. We set the marginal variance of the measurement error $\zeta_1 \in \{0.005, 0.01\}$. ARL_{target} is set to 10,000 which is approximately equivalent to 100 days as the inter-reporting time is set to 15 minutes in the SWMM model. The S^3T statistic defined in (57) with a window length $\omega = 20$ is used as the detection statistic. The search space of the parameter θ in a VAR(1) model is set to $\{0.1, 0.2, \dots, 0.9\}$. For CSCBO, we use the same setting as in Section 4.4.3.1. The minimum requirement for reliability is set as $q = 0.9$.

4.4.3.2 Comparison on detection delay

To demonstrate the advantage of the S^3T statistic, we compare S^3T with two other detection statistics in terms of conditional expected detection delay for fixed sets of sensors locations when the target ARL_0 is set to a pre-specified level: (i) a simple detection statistic defined in (52) which we refer to as the Shewhart statistic and (ii) a max-CUSUM statistic defined as follows,

$$\text{Cusum}_{t,i} = \max(0, \text{Cusum}_{t-1,i} + c_i - \kappa\sqrt{\zeta_1}), \quad i = 1, \dots, M, \quad (60)$$

and

$$W_t = \max_{i \in \{1, \dots, M\}} \text{Cusum}_{t,i}, \quad (61)$$

Table 12: Simulated ARL_0 using S^3T when no contaminant event. $ARL_{\text{target}} = 10000$.

M	ζ_1	b	ARL_0
5	0.005	4.989	10276
5	0.01	5.025	10426
7	0.005	4.788	9904
7	0.01	4.847	10143

where $\kappa > 0$ is a constant and we set $\kappa = 0.1$ in the experiment. Here we assume that the marginal variance of the sensor measurement error ζ_1 is known in the max-CUSUM statistic, while in practice, it should be estimated from data. The threshold b for the three statistics are adjusted via simulations so that the ARLs are approximately 10,000.

We use the two sensor placements marked by circles in Figure 18(a) ($M = 5$) and (c) ($M = 7$) in the experiments. Note that these solutions are not necessarily feasible for problem (56) and the focus of this experiment is on the performance of the S^3T statistic alone. We compare the ECEDDs of each sensor placement using different detection statistics. The ECEDDs and the corresponding standard error are obtained based on 1000 SWMM runs and presented in Table 11. As we can observe, S^3T outperforms the other two statistics in all experiments and achieves 7 - 12 hours shorter conditional expected detection delay comparing to the Shewhart statistic and 1 - 3 hours comparing to the max-CUSUM statistic. The advantage is more obvious when the variance of measurement error is high ($\zeta_1 = 0.01$). Such an advantage of S^3T is due to its capacity of capturing both spatial and temporal correlation in the data.

4.4.3.3 Adjusting threshold b

Here we validate the accuracy of the simulation approach for adjusting b discussed in Section 4.3.2. We test the two sensor placements marked by circles in Figure 18(a) ($M = 5$) and (c) ($M = 7$). In the experiment, we first adjust b using the discussed approach with $ARL_{\text{target}} = 10000$ and then simulate in-control data to obtain the actual ARL_0 . The number of replications for estimating ARL_0 is 2000. The results is presented in Table 12 where we can observe that the adjusted b values are accurate enough for the detection

Table 13: Performance metrics of the optimal feasible solution found by CSCBO with S^3T as the detection statistic. The unit of ECEDD is hour.

M	ζ_1	ECEDD	ER	false_alarm_rate	NUM
5	0.01	70.60	0.931	0.2%	147
5	0.005	64.36	0.924	0.26%	181
7	0.01	62.46	0.934	0.32%	123
7	0.005	57.93	0.925	0.31%	253

statistic to achieve the target ARL_0 .

4.4.3.4 Optimal solution in the presence of measurement error

Now we combine CSCBO and S^3T to solve problem (56). In each iteration of the optimization, the threshold b is adjusted so that the ARL_0 is approximately 10,000. Figure 19 presents the optimal feasible solutions found by CSCBO with S^3T as the detection statistics for $M = 5$ and $M = 7$. Table 13 reports the performance metrics of these solutions based on 1000 SWMM runs. Note that in the experiments, a false alarm is defined as an alarm raised by the sensors before the actual starting time of a simulated contaminant event. We calculate false alarm rates as the number of false alarms divided by the total number of simulated contaminant events. From Table 13, we first observe that all solutions have estimated reliability higher than $q = 0.9$, and hence the reported solutions are feasible. The false alarm rates are around 0.3% in all cases, indicating that by controlling the ARL_0 of sensor networks, the occurrence of false alarms can be controlled at a very low level. In addition, we can observe from Figure 19 that the optimal sensor placements found by the algorithm do not differ much for different levels of the variance of measurement error. This means that the optimal sensor network design identified by the combined procedure of CSCBO and S^3T is robust to different level of variances of the sensor measurement error.

4.5 Conclusion

In this chapter, we study the problem of optimal sensor network design in the presence of random measurement error for water quality monitoring on river systems. The spatial and temporal correlation in the sensor measurement is take into consideration. We formulate

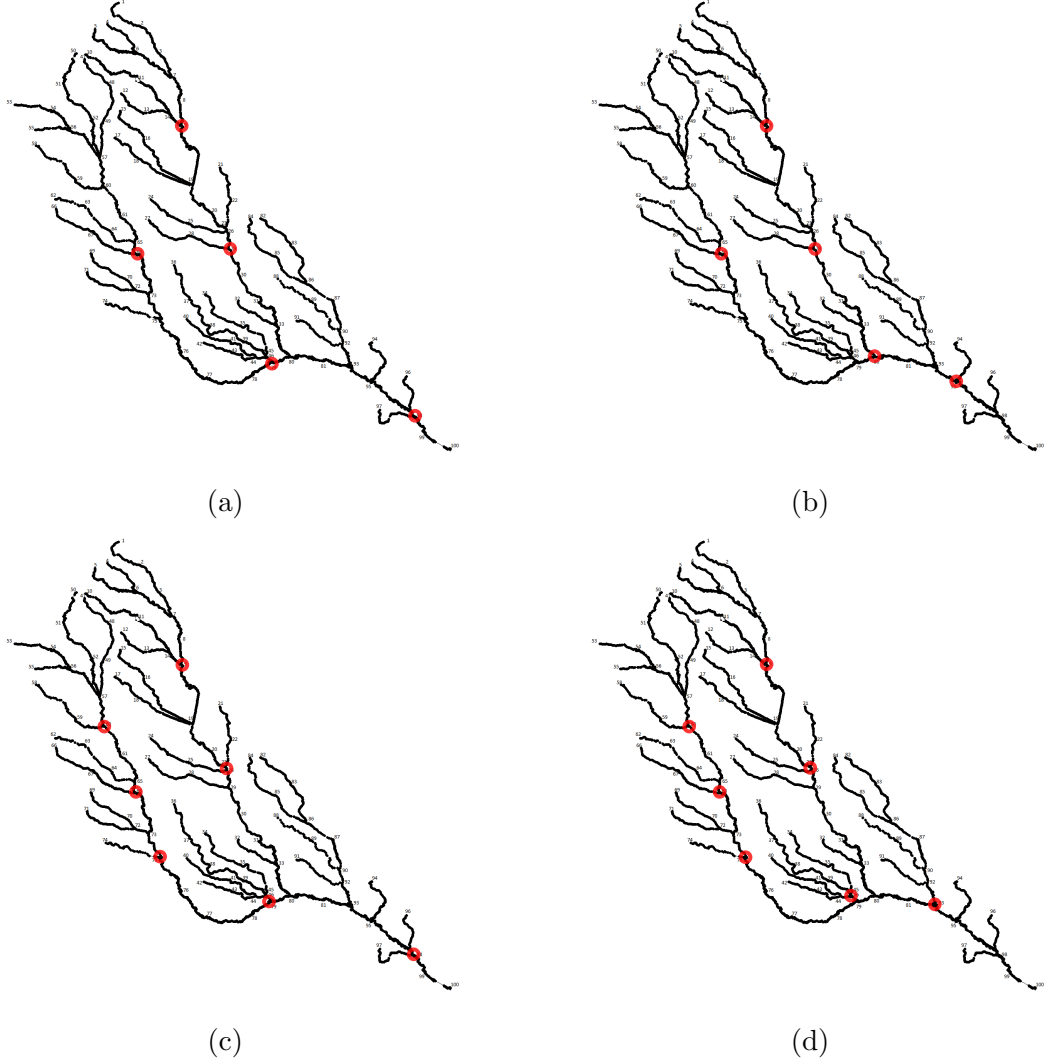


Figure 19: Optimal feasible solutions found by CSCBO with S^3T as the detection statistic: (a) $M = 5$ and $\zeta_1 = 0.01$, (b) $M = 5$ and $\zeta_1 = 0.005$, (c) $M = 7$ and $\zeta_1 = 0.01$ and (d) $M = 7$ and $\zeta_1 = 0.005$.

the problem as a joint problem of constrained black-box function optimization and spatio-temporal change-point detection. The proposed algorithm CSCBO with a Wasserstein similarity metric is demonstrated to converge quickly and is able to find competitive feasible suboptimal sensor placements. By combining CSCBO with the S^3T statistic, we identify sensor network designs which are shown to have low conditional detection delay, required reliability and low false-alarm rates. In addition, the identified sensor networks are robust to different levels of measurement error variances.

CHAPTER V

DISCUSSION AND CONCLUSIONS

In Chapter 2, we discuss dimension reduction via spatial scanning. The key question we aim to answer is how much do we lose in terms of detection performance by using reduced-dimension spatial scanning comparing to methods using full observation vectors. Our studies show that the performance loss of the RD approach is within the acceptable range. Considering that the RD spatial scanning enjoys all the computation and implementation benefits, it should be a preferable method for practitioners.

In Chapter 3, we propose the S^3T statistic to detect the emergence of a spatially and temporally correlated signal from noisy background. The proposed statistic jointly captures the spatial and temporal correlations of the signal. Numerical studies based on simulated and real data demonstrate that S^3T outperforms the baseline methods. Our results also show that the advantage of S^3T is more obvious when the signal to be detected is weak.

In Chapter 4, we study the problem of optimal sensor network design, which is formulated as a joint problem of constrained black-box function optimization and spatio-temporal change-point detection. We propose the Confidence-Set based Constrained Bayesian Optimization (CSCBO) algorithm for general constrained black-box function optimization problems and extend the algorithm to tackle with high-dimensional binary decision variables by using the Wasserstein similarity metric. Experiments reveal that confidence-set is an effective and flexible method to handle black-box function constraints. In addition, CSCBO converges much quicker than NP+PFM in the water monitoring network application and can identify competitive local optimal solutions with little loss on the conditional expected detection delay. Finally, sensor network designs identified by combining CSCBO and S^3T are shown to be robust to sensor measurement errors.

APPENDIX A

APPENDIX FOR CHAPTER 2

A.1 Derivation of (12)

First note that $d_0 < 0$ in all cases and we assume that $d_{c,r} > 0$. Given target ARL_0 , d_0 , $d_{c,r}$, Ω_0^2 and $\Omega_{c,r}^2$, we need to find H and ARL_1 that satisfy the following approximations from (11):

$$\text{ARL}_0 \approx \frac{\Omega_0^2}{2d_0^2} \left\{ \exp \left[-\frac{2d_0(H + 1.166\Omega_0)}{\Omega_0^2} \right] - 1 + \frac{2d_0(H + 1.166\Omega_0)}{\Omega_0^2} \right\};$$

and

$$\text{ARL}_1 \approx \frac{\Omega_{c,r}^2}{2d_{c,r}^2} \left\{ \exp \left[-\frac{2d_{c,r}(H + 1.166\Omega_{c,r})}{\Omega_{c,r}^2} \right] - 1 + \frac{2d_{c,r}(H + 1.166\Omega_{c,r})}{\Omega_{c,r}^2} \right\}.$$

From the above two approximations, we get

$$\frac{2d_0^2}{\Omega_0^2} \text{ARL}_0 + 1 \approx \left\{ \exp \left[-\frac{2d_0(H + 1.166\Omega_0)}{\Omega_0^2} \right] + \frac{2d_0(H + 1.166\Omega_0)}{\Omega_0^2} \right\}; \quad (62)$$

and

$$\frac{2d_{c,r}^2}{\Omega_{c,r}^2} \text{ARL}_1 + 1 \approx \left\{ \exp \left[-\frac{2d_{c,r}(H + 1.166\Omega_{c,r})}{\Omega_{c,r}^2} \right] + \frac{2d_{c,r}(H + 1.166\Omega_{c,r})}{\Omega_{c,r}^2} \right\}. \quad (63)$$

Let $\eta_0 = \frac{2d_0^2}{\Omega_0^2} \text{ARL}_0 + 1$ and $\eta_1 = \frac{2d_{c,r}^2}{\Omega_{c,r}^2} \text{ARL}_1 + 1$. In general, a solution x to an equation $e^x + x = c$ for a constant c can be expressed using the Lambert W function [10]. Then using the Lambert W function, we get

$$\frac{2d_0(H + 1.166\Omega_0)}{\Omega_0^2} = W_{-1}(-e^{-\eta_0}) + \eta_0 \quad \text{from (62),} \quad (64)$$

$$\frac{2d_{c,r}(H + 1.166\Omega_{c,r})}{\Omega_{c,r}^2} = W_0(-e^{-\eta_1}) + \eta_1 \quad \text{from (63),} \quad (65)$$

where $W_{-1}(\cdot)$ and $W_0(\cdot)$ denote the two branches of the Lambert W function [10]. Note that we need $W_{-1}(\cdot)$ for (64) because the LHS is a negative number due to $d_0 < 0$ while we need $W_0(\cdot)$ for (65) because the LHS is a positive number due to $d_{c,r} > 0$.

Solving (64) and (65) for H results in

$$H = \frac{\Omega_0^2}{2d_0}(W_{-1}(-e^{-\eta_0}) + \eta_0) - 1.166\Omega_0$$

and

$$H = \frac{\Omega_{c,r}^2}{2d_{c,r}}(W_0(-e^{-\eta_1}) + \eta_1) - 1.166\Omega_{c,r}.$$

Then we obtain

$$\frac{\Omega_0^2}{2d_0}(W_{-1}(-e^{-\eta_0}) + \eta_0) - 1.166\Omega_0 = \frac{\Omega_{c,r}^2}{2d_{c,r}}(W_0(-e^{-\eta_1}) + \eta_1) - 1.166\Omega_{c,r}.$$

Let $\epsilon_{\eta_0} = -W_{-1}(-e^{-\eta_0}) - \eta_0$, which is a positive constant number. Then

$$\frac{\Omega_{c,r}^2}{2d_{c,r}}(W_0(-e^{-\eta_1}) + \eta_1) = -\frac{\Omega_0^2}{2d_0}\epsilon_{\eta_0} - 1.166(\Omega_0 - \Omega_{c,r}). \quad (66)$$

The LHS of (66) is

$$\frac{\Omega_{c,r}^2}{2d_{c,r}} \left(W_0(-e^{-\eta_1}) + \frac{2d_{c,r}^2}{\Omega_{c,r}^2} \text{ARL}_1 + 1 \right).$$

Note that $-1 \leq W_0(-e^{-\eta_1}) \leq 0$ and $1 + W_0(-e^{-\eta_1}) \approx 0$. Then we finally obtain (12) as follows,

$$\text{ARL}_1 \approx -\frac{\Omega_0^2}{2d_0d_{c,r}}\epsilon_{\eta_0} - 1.166\frac{1}{d_{c,r}}(\Omega_0 - \Omega_{c,r}).$$

A.2 Figures and Results for Section 2.3.2.1

Numerical Examples with Known Center and Known Radius: Here we conduct simulation experiments assuming the shift cluster is known and show that simulation results match the ARL_1 measures very well.

Consider two different spatial covariance matrices: (i) tridiagonal matrix $\Sigma_1(\rho)$ and (ii) $\Sigma_2(\rho) \in R_{p \times p}$ with $[\Sigma_2(\rho)]_{i,j} = \rho^{|i-j|}$, $\forall i, j \in \{1, \dots, p\}$, where the correlation decays polynomially. We use $p = 49$, $\tilde{p} = 5$ and $\rho \in \{0, 0.02, 0.04, \dots, 0.3\}$. For the out-of-control state, homogeneous shifts (shifts of all affected locations in the cluster have same magnitude) with magnitudes $\delta = 0.25, 0.5, 0.75$ and 1 are tested. The targeted ARL_0 is fixed to 1000 in all the cases.

We report the ratio of the simulated ARL_1 of a F-MCUSUM chart to a RD-MCUSUM chart. The result is shown in Figure 20. In the same plot, we also present the ratio of the

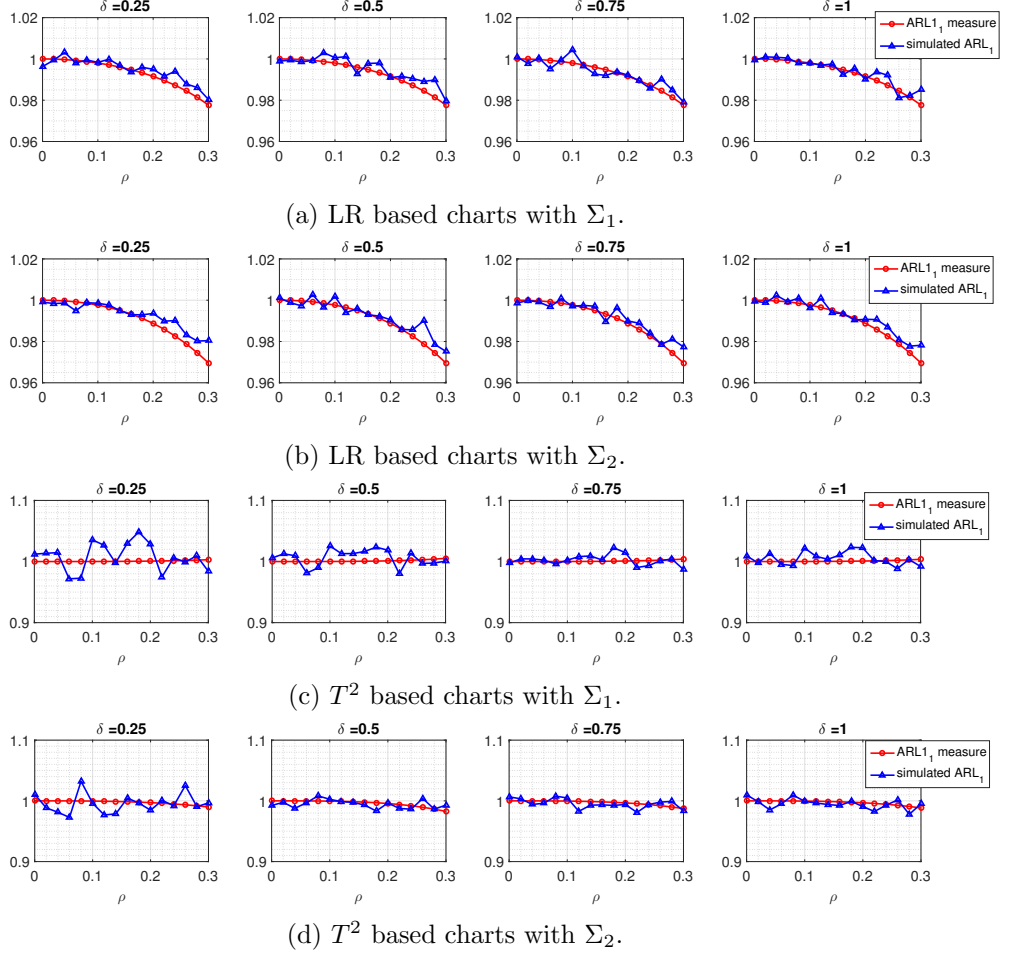


Figure 20: Ratio of simulated ARL_1 's (blue) and ratio of ARL_1 measures (red) in the known shift cluster case.

ARL_1 measures, m_{LR}/\tilde{m}_{LR} and m_{T^2}/\tilde{m}_{T^2} . Note that for T^2 based charts, when the shift magnitude $\delta = 0.25$, the out-of-control drift $d_{c,r}$ is negative. Hence, the ARL_1 measures are not applicable in this case. Instead, we numerically search for the control limit H given a ARL_0 , obtain d_0 and Ω_0 using Formula (11), and then substitute H into (11) with $d_{c,r}$ and $\Omega_{c,r}$ to obtain an approximate ARL_1 for both T^2 -F-MCUSUM and T^2 -RD-MCUSUM charts.

From Figure 20, we see that simulation results match very well with the ARL_1 measures. This indicates that the ARL_1 measure is indeed a good measure for detection performance comparison. Moreover, we observe that the F-MCUSUM and RD-MCUSUM charts have similar detection performances in the known shift cluster case: (i) for LR based methods,

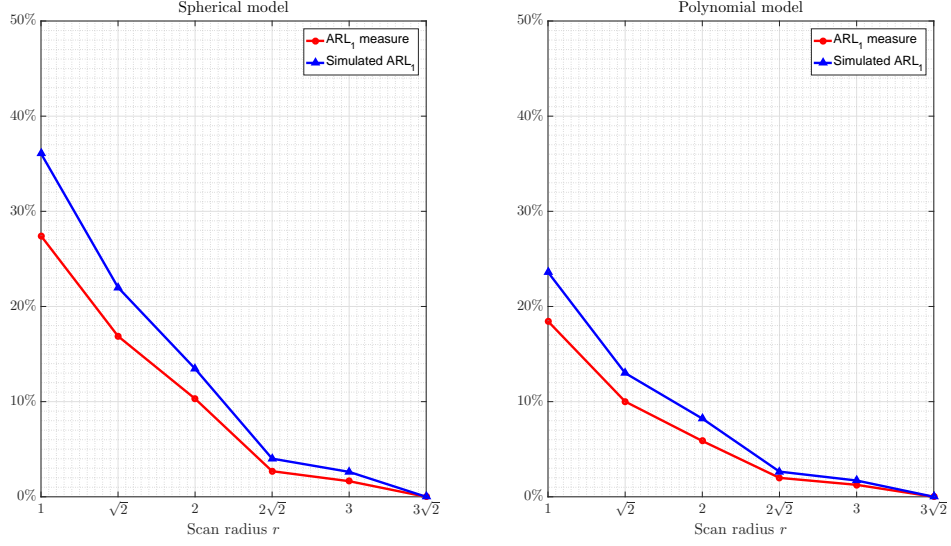


Figure 21: Performance loss based on ARL_1 measure and simulated ARL_1 .

the F-MCUSUM chart slightly outperform RD-MCUSUM chart since all the ratios is less than but very close to 1; (ii) for T^2 based methods, depending on different spatial correlation structures, m_{T^2}/\tilde{m}_{T^2} is either slightly greater or smaller than 1, however, we can hardly see the difference in performance from simulation results.

A.3 Decaying Shift

Here we consider “bell” shaped shift signals. Such a shift signal occurs at a center $c \in P$, affect all locations or sensors, and the shift magnitude decays with the distances from c . In this case, reduced-dimension charts are destined to lose performance, since each reduced-dimension observation can only capture a portion of the signal energy. The smaller the scanning radius, the more dimensionality reduction we achieve since each cluster has fewer data streams, but there is a great loss of performance.

We present a simple example in the presence of such a decaying shift. This example demonstrates one practical use of our theory: analytically choose the smallest scan radius r allowed based on the ARL_1 measure, given a maximum acceptable performance loss. Here we use percentage loss as the metric:

$$\frac{ARL_1 \text{ of an RD-MCUSUM chart} - ARL_1 \text{ of a F-MCUSUM chart}}{ARL_1 \text{ of an F-MCUSUM chart}} \times 100\%,$$

and compare it with the ratio with approximated ARL_1 measures defined as $(\tilde{m}_{LR} - m_{LR})/m_{LR}$ using (13). Consider a shift signal with a magnitude δ at the center c , and the magnitude at other locations $d \in P$ is given by $\delta\theta^{\|q_c - q_d\|}$, where $\theta \in (0, 1)$ is the delay rate. Assume a monitoring area of dimension $p = 7 \times 7$. We run the control charts on the spherical model and the polynomial model with correlation parameter $\rho = 0.2$. We let $\delta = 1$, $\theta = 0.5$ and target $\text{ARL}_0 = 1000$. Our goal is to determine a smallest scan radius r (i.e., maximum dimensionality reduction), such that the performance loss in ARL_1 is no more than, say, 10% when using reduced dimension vectors. We may choose r from $\{0, 1, \sqrt{2}, 2, 2\sqrt{2}, 3, 3\sqrt{2}\}$. Note that $r = 3\sqrt{2}$ corresponds to using full observation vectors, and $r = 0$ means each scanning region has only one sensor. In Figure 21, the blue curve shows the performance loss in terms of the simulated ARL_1 , while the red curve shows $(\tilde{m}_{LR} - m_{LR})/m_{LR}$. From Figure 21, we conclude that we need $r = 2$ or $2\sqrt{2}$ under Σ_3 and $r = \sqrt{2}$ or 2 under Σ_4 to achieve no more than 10% loss in ARL_1 .

APPENDIX B

APPENDIX FOR CHAPTER 3

B.1 Derivation of $\frac{\partial \ell}{\partial \mu} \Big|_{\mu=0, \gamma=0}$.

The following propositions are used in the derivation of $\frac{\partial \ell}{\partial \mu} \Big|_{\mu=0, \gamma=0}$.

Proposition B.1.1. *Let $\mathbf{M}(t)$ be a nonsingular square matrix whose elements are functions of a scalar parameter α . Then,*

$$\frac{\partial \mathbf{M}(t)^{-1}}{\partial \alpha} = -\mathbf{M}(t)^{-1} \frac{\partial \mathbf{M}(t)}{\partial \alpha} \mathbf{M}(t)^{-1}.$$

Proposition B.1.2. *Let $\mathbf{M}(t)$ be a nonsingular square matrix whose elements are functions of a scalar parameter α . Then,*

$$\frac{\partial |\mathbf{M}(t)|}{\partial \alpha} = |\mathbf{M}(t)| \text{tr} \left(\mathbf{M}(t)^{-1} \frac{\partial \mathbf{M}(t)}{\partial \alpha} \right).$$

By Proposition B.1.1, we can calculate,

$$\begin{aligned} \left. \frac{\log |\gamma \mathbf{V}_\tau(\theta) + \boldsymbol{\Sigma}_\tau|}{\partial \gamma} \right|_{\mu=0, \gamma=0} &= \frac{1}{|\gamma \mathbf{V}_\tau(\theta) + \boldsymbol{\Sigma}_\tau|} \left. |\gamma \mathbf{V}_\tau(\theta) + \boldsymbol{\Sigma}_\tau| \text{tr} \left((\gamma \mathbf{V}_\tau(\theta) + \boldsymbol{\Sigma}_\tau)^{-1} \mathbf{V}_\tau(\theta) \right) \right|_{\gamma=0} \\ &= \text{tr}(\boldsymbol{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta)). \end{aligned}$$

For convenience, here we use \mathbf{y} and $\boldsymbol{\mu}$ to denote $\mathbf{y}_{(k+1:N)}$ and $\boldsymbol{\mu}_{(k+1:N)}$. By Proposition B.1.2, we have,

$$\begin{aligned} \left. \frac{\partial (\mathbf{y} - \boldsymbol{\mu})^\top (\gamma \mathbf{V}_\tau(\theta) + \boldsymbol{\Sigma}_\tau)^{-1} (\mathbf{y} - \boldsymbol{\mu})}{\partial \gamma} \right|_{\mu=0, \gamma=0} &= (\mathbf{y} - \boldsymbol{\mu})^\top \frac{\partial (\gamma \mathbf{V}_\tau(\theta) + \boldsymbol{\Sigma}_\tau)^{-1}}{\partial \gamma} (\mathbf{y} - \boldsymbol{\mu}) \Big|_{\mu=0, \gamma=0} \\ &= -\mathbf{y}^\top (\gamma \mathbf{V}_\tau(\theta) + \boldsymbol{\Sigma}_\tau)^{-1} \mathbf{V}_\tau(\theta) (\gamma \mathbf{V}_\tau(\theta) + \boldsymbol{\Sigma}_\tau)^{-1} \mathbf{y} \Big|_{\gamma=0} \\ &= -\mathbf{y}^\top \boldsymbol{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta) \boldsymbol{\Sigma}_\tau^{-1} \mathbf{y}. \end{aligned}$$

Hence we have,

$$\frac{\partial \ell}{\partial \mu} \Big|_{\mu=0, \gamma=0} = -\frac{1}{2} \text{tr}(\boldsymbol{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta)) + \frac{1}{2} \mathbf{y}^\top \boldsymbol{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta) \boldsymbol{\Sigma}_\tau^{-1} \mathbf{y},$$

as appeared in equation (24).

B.2 Derivation of $\text{Var}[S(\tau, \theta)]$.

Here we calculate the variance of the statistic $S(\tau, \theta)$ defined in (25). For convenience, we use \mathbf{y} to denote $\mathbf{y}_{(k+1:N)}$, use \mathbf{V} to denote the matrix $\Sigma_\tau^{-1} \mathbf{V}_\tau(\theta) \Sigma_\tau^{-1}$, use Σ to denote Σ_τ and use c and d to denote $c(\tau, \theta)$ and $d(\tau, \theta)$, respectively. Then we can write,

$$S(\tau, \theta) = \frac{(\mathbf{y}^\top \mathbf{V} \mathbf{y} - c)^2}{d} + \mathbf{y}^\top \Sigma^{-1} \mathbf{y}.$$

We have that $\text{E}[S] = p\tau + 1$, and

$$\text{Var}[S(\tau, \theta)] = \text{E}[S^2] - \text{E}[S]^2.$$

In the following, we calculate $\text{E}[S^2]$.

$$\begin{aligned} \text{E}[S^2] &= \text{E} \left[\left(\frac{(\mathbf{y}^\top \mathbf{V} \mathbf{y} - c)^2}{d} + \mathbf{y}^\top \Sigma^{-1} \mathbf{y} \right)^2 \right] \\ &= \text{E} \left[\left(\mathbf{y}^\top \Sigma^{-1} \mathbf{y} \right)^2 \right] + 2\text{E} \left[\frac{(\mathbf{y}^\top \mathbf{V} \mathbf{y} - c)^2}{d} \mathbf{y}^\top \Sigma^{-1} \mathbf{y} \right] + \text{E} \left[\frac{(\mathbf{y}^\top \mathbf{V} \mathbf{y} - c)^4}{d^2} \right]. \end{aligned} \quad (67)$$

The first term can be calculated as,

$$\text{E} \left[\left(\mathbf{y}^\top \Sigma^{-1} \mathbf{y} \right)^2 \right] = p^2 \tau^2 + 2p\tau. \quad (68)$$

We then expand the second term,

$$\begin{aligned} \text{E} \left[\frac{(\mathbf{y}^\top \mathbf{V} \mathbf{y} - c)^2}{d} \mathbf{y}^\top \Sigma^{-1} \mathbf{y} \right] &= \frac{1}{d} \text{E} \left[\left(\mathbf{y}^\top \mathbf{V} \mathbf{y} \right)^2 \mathbf{y}^\top \Sigma^{-1} \mathbf{y} \right] \\ &\quad - \frac{2c}{d} \text{E} \left[\mathbf{y}^\top \mathbf{V} \mathbf{y} \mathbf{y}^\top \Sigma^{-1} \mathbf{y} \right] + \frac{c^2}{d} \text{E} \left[\mathbf{y}^\top \Sigma^{-1} \mathbf{y} \right]. \end{aligned}$$

We calculate the three expectations separately:

$$\text{E} \left[\mathbf{y}^\top \Sigma^{-1} \mathbf{y} \right] = p\tau.$$

$$\text{E} \left[\mathbf{y}^\top \mathbf{V} \mathbf{y} \mathbf{y}^\top \Sigma^{-1} \mathbf{y} \right] = (p\tau + 2)c.$$

$$\text{E} \left[\left(\mathbf{y}^\top \mathbf{V} \mathbf{y} \right)^2 \mathbf{y}^\top \Sigma^{-1} \mathbf{y} \right] = (p\tau + 4)(2d + c^2).$$

Combining we get,

$$\text{E} \left[\frac{(\mathbf{y}^\top \mathbf{V} \mathbf{y} - c)^2}{d} \mathbf{y}^\top \Sigma^{-1} \mathbf{y} \right] = 2p\tau + 4. \quad (69)$$

Next we calculate the last term in (67),

$$\mathbb{E}\left[\frac{(\mathbf{y}^\top \mathbf{V} \mathbf{y} - c)^4}{d^2}\right] = 3 - 2p\tau - 24\frac{c}{d^2}\text{tr}\left(\Sigma^{-1}\mathbf{V}\Sigma^{-1}\mathbf{V}\Sigma^{-1}\mathbf{V}\right) + \frac{48}{d^2}\text{tr}\left(\Sigma_\tau^{-1}\mathbf{V}\Sigma^{-1}\mathbf{V}\Sigma^{-1}\mathbf{V}\Sigma^{-1}\mathbf{V}\right). \quad (70)$$

Note that the tedious calculation steps for (70) are omitted here.

Combining (68), (69) and (70), we can obtain,

$$\text{Var}[S(\tau, \theta)] = 2p\tau + 10 - 24\frac{c}{d^2}\text{tr}\left(\Sigma^{-1}\mathbf{V}\Sigma^{-1}\mathbf{V}\Sigma^{-1}\mathbf{V}\right) + \frac{48}{d^2}\text{tr}\left(\Sigma^{-1}\mathbf{V}\Sigma^{-1}\mathbf{V}\Sigma^{-1}\mathbf{V}\Sigma^{-1}\mathbf{V}\right).$$

B.3 Derivation of the cumulant generating function of W .

Here we present the derivation of the cumulant generating function of $W(\tau, \theta)$ under the null hypothesis, i.e. equation (34).

Let $\mathbf{z} = \Sigma_\tau^{-\frac{1}{2}} \mathbf{y}_{(k+1:N)}$. Under the null hypothesis, we have $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{p\tau})$. For convenience, here we use \mathbf{B} to denote the $p\tau$ by $p\tau$ matrix $\Sigma_\tau^{-\frac{1}{2}} \mathbf{V}_\tau(\theta) \Sigma_\tau^{-\frac{1}{2}}$, and use c and d to denote $c(\tau, \theta)$ and $d(\tau, \theta)$, respectively. Then, we have

$$W(\tau, \theta) = \frac{\mathbf{z}^\top \mathbf{B} \mathbf{z} - c}{\sqrt{d}}.$$

Under the null hypothesis, the cumulant generating function of $W(\tau, \theta)$ can be calculated as

$$\begin{aligned} \psi(\xi) &= \log \mathbb{E}[\exp(\xi W(\tau, \theta))] = \log \mathbb{E}\left[\exp\left(\xi \left(\frac{\mathbf{z}^\top \mathbf{B} \mathbf{z} - c}{\sqrt{d}}\right)\right)\right] \\ &= -\xi \frac{c}{\sqrt{d}} + \log \mathbb{E}\left[\exp\left(\frac{\xi \mathbf{z}^\top \mathbf{B} \mathbf{z}}{\sqrt{d}}\right)\right] \\ &= -\xi \frac{c}{\sqrt{d}} + \log \int_{\mathbf{z}} \exp\left(\frac{\xi \mathbf{z}^\top \mathbf{B} \mathbf{z}}{\sqrt{d}}\right) \frac{1}{(2\pi)^{\frac{p\tau}{2}}} \exp\left(-\frac{1}{2} \mathbf{z}^\top \mathbf{z}\right) d\mathbf{z} \\ &= -\xi \frac{c}{\sqrt{d}} + \log \int_{\mathbf{z}} \frac{1}{(2\pi)^{\frac{p\tau}{2}}} \exp\left(-\frac{1}{2} \mathbf{z}^\top \left(\mathbf{I}_{p\tau} - \frac{2\xi \mathbf{B}}{\sqrt{d}}\right) \mathbf{z}\right) d\mathbf{z} \\ &= -\xi \frac{c}{\sqrt{d}} + \log \left| \mathbf{I}_{p\tau} - \frac{2\xi \mathbf{B}}{\sqrt{d}} \right|^{-\frac{1}{2}}, \end{aligned}$$

which is equivalent to equation (34). Note that the last equation uses the fact that

$$\int_{\mathbf{z}} \frac{1}{(2\pi)^{\frac{p\tau}{2}}} \exp\left(-\frac{1}{2} \mathbf{z}^\top \left(\mathbf{I}_{p\tau} - \frac{2\xi \mathbf{B}}{\sqrt{d}}\right) \mathbf{z}\right) d\mathbf{z} = \left| \mathbf{I}_{p\tau} - \frac{2\xi \mathbf{B}}{\sqrt{d}} \right|^{-\frac{1}{2}}.$$

B.4 Proof of Theorem 3.3.1

After discretizing the parameter space, $W(\tau, \theta)$ is treated as a two-dimensional Gaussian random field, which is completely characterized by its covariance function. The following

lemma computes the covariance function of $W(\tau, \theta)$.

Lemma B.4.1. *Under the null hypothesis, the covariance function of $W(\tau, \theta)$ is*

$$\text{Cov}[W(n, \theta_1), W(m, \theta_2)] = \frac{\text{tr}(\mathbf{A}_n(\theta_1)\mathbf{A}_n(\theta_2))}{\left[\text{tr}(\mathbf{A}_n(\theta_1)\mathbf{A}_n(\theta_1))\text{tr}(\mathbf{A}_m(\theta_2)\mathbf{A}_m(\theta_2))\right]^{1/2}}, \quad (71)$$

where $n \leq m$.

The following lemma shows that the first order approximation of the covariance function in (71) does not have any cross product term. Thus, the two-dimensional random field is further decomposed as a sum of two independent one-dimensional random processes.

Lemma B.4.2. *Assuming that δ and $i \in Z$ are small relative to θ and τ , respectively, the first order approximation of the covariance function in (71) is given as,*

$$\text{Cov}[W(\tau, \theta), W(\tau + i, \theta + \delta)] \approx 1 - \gamma^2(\tau, \theta)\delta^2 - \frac{\mu(\tau, \theta)}{2\tau}i + o(\delta^2) + o(i), \quad (72)$$

where

$$\gamma(\tau, \theta) = \frac{\text{tr}(\dot{\mathbf{A}}_\tau(\theta)\mathbf{A}_\tau(\theta))}{\text{tr}(\mathbf{A}_\tau(\theta)\mathbf{A}_\tau(\theta))}, \quad (73)$$

$\mu(\tau, \theta)$ is defined in (32), and $\dot{\mathbf{A}}_\tau(\theta) = \partial\mathbf{A}_\tau(\theta)/\partial\theta$.

The following two Lemmas are needed in the proof. Both Lemmas are proved in [73].

Lemma B.4.3. *Assume $\xi \rightarrow \infty$, $b \rightarrow \infty$, $N \rightarrow \infty$, with $\frac{\xi}{b} \approx 1$ and $\frac{b}{N} \approx c$, where $c > 0$ is some constant. The discretized process $b\left[W\left(\tau + i, \theta + \frac{\Delta}{\sqrt{N}j}\right) - \xi\right]$, where i is an integer and $j \geq 0$, conditioned on $W(\tau, \theta) = \xi$ can be written as a sum of two independent processes:*

$$\left\{b\left[W\left(\tau + i, \theta + \frac{\Delta}{\sqrt{N}j}\right) - \xi\right] \middle| W(\tau, \theta) = \xi\right\} = S_i + V_j,$$

where $S_i = \sum_{l=1}^i a_\ell$ with

$$a_\ell \sim N\left(-\frac{\mu(\tau, \theta)}{2\tau}b^2, \frac{\mu(\tau, \theta)}{\tau}b^2\right),$$

and

$$V_j = \sqrt{2}\gamma(\tau, \theta)\frac{b}{\sqrt{N}}\Delta jV - \gamma^2(\tau, \theta)\frac{b^2}{N}\Delta^2 j^2,$$

with $V \sim N(0, 1)$. $\mu(\tau, \theta)$ and $\gamma(\tau, \theta)$ are defined in (32) and (73), respectively.

Lemma B.4.4. Assume x_1, x_2, \dots are i.i.d. $N(-\mu_1, \sigma_1^2)$ random variables ($\mu_1 > 0$). Define the random walk $S_0 = 0, S_i = \sum_{l=1}^i x_l, i = 1, 2, \dots$, and the smooth varying random process $V_j = \beta \Delta j V - \frac{\beta^2}{2} \Delta^2 j^2$, for some constants $\Delta > 0, \beta > 0$. As $\Delta \rightarrow 0$, for some constant α , we have

$$\frac{1}{\Delta} \int_0^\infty e^{-\alpha x} \mathbb{P}\left(\max_{i \geq 1} S_i \leq -x\right) \mathbb{P}\left(\max_{i \leq 0} S_i + \max_{j \geq 1} V_j \leq -x\right) dx \xrightarrow{\Delta \rightarrow 0} \frac{|\beta|}{\sqrt{2\pi}} \left(\frac{2\mu_1^2}{\sigma_1^2}\right) \nu\left(\frac{2\mu_1}{\sigma_1}\right),$$

where $\nu(x)$ is defined in (31).

In the following, we go through the main steps that lead to the approximation of the false alarm rate in Theorem 3.3.1 for the case of $d = 1$.

Step 1: We first discretize the parameter $\theta \in [\theta_1, \theta_2]$ by a rectangular mesh grid of size $\frac{\Delta}{\sqrt{N}}$, where $\Delta > 0$ is a small number. Note that the discretization mentioned here is used for asymptotic analysis only. The probability of false alarm can be approximated as

$$\mathbb{P}\left(\max_{(i,j) \in D} W\left(i, j \frac{\Delta}{\sqrt{N}}\right) \geq b\right), \quad (74)$$

where D is the index set

$$D = \left\{(i, j) : 1 \leq i \leq N, \theta_1 \leq j \frac{\Delta}{\sqrt{N}} \leq \theta_2\right\},$$

which covers the entire parameter space. Let $J(i_0, j_0)$ denote everything to the “future” of the current index (i_0, j_0) in the parameter space, i.e.,

$$J(i_0, j_0) = \{(i, j) \in D : j \geq j_0, \text{ or } i \geq i_0 \text{ and } j = j_0\}.$$

Using the similar approach as in [54], the event

$$\left\{\max_{(i,j) \in D} W\left(i, j \frac{\Delta}{\sqrt{N}}\right) \geq b\right\}$$

can be decomposed into a series of “last hitting events”, for which (i_0, j_0) is the “last” location where $W\left(i, j \frac{\Delta}{\sqrt{N}}\right)$ hits the threshold b . Then, the probability in (74) can be written as the sum of probabilities of $W\left(i, j \frac{\Delta}{\sqrt{N}}\right)$ last hits b at (i_0, j_0) over all possible

(i_0, j_0) :

$$\begin{aligned}
\mathbb{P}\left(\max_{(i,j) \in D} W\left(i, j \frac{\Delta}{\sqrt{N}}\right) \geq b\right) &\approx \sum_{(i_0, j_0) \in D} \mathbb{P}\left(W\left(i_0, j_0 \frac{\Delta}{\sqrt{N}}\right) \geq b, \max_{(i,j) \in J(i_0, j_0)} W\left(i, j \frac{\Delta}{\sqrt{N}}\right) < b\right) \\
&= \sum_{(i_0, j_0) \in D} \int_0^\infty \mathbb{P}\left(W\left(i_0, j_0 \frac{\Delta}{\sqrt{N}}\right) = b + \frac{x}{b}\right) \\
&\quad \times \mathbb{P}\left(\max_{(i,j) \in J(i_0, j_0)} W\left(i, j \frac{\Delta}{\sqrt{N}}\right) < b \mid W\left(i_0, j_0 \frac{\Delta}{\sqrt{N}}\right) = b + \frac{x}{b}\right) \frac{dx}{b}.
\end{aligned} \tag{75}$$

Step 2: In the following, we obtain an approximation on the probability

$$\mathbb{P}\left(W\left(i_0, j_0 \frac{\Delta}{\sqrt{N}}\right) = b + \frac{x}{b}\right) \frac{dx}{b}.$$

To simplify the notation, we denote $W\left(i_0, j_0 \frac{\Delta}{\sqrt{N}}\right)$ as W here. The key idea is to approximate W as a Gaussian random field. The Gaussian approximation performs well when the probability of interest is close to the mean of the true distribution, but suffers from deviation if the probability is in the tail of the true distribution. Hence, we apply the change-of-measure technique to shift the mean of the random field W to the threshold b .

Denote the cumulant generating function of W as $\psi(\xi) = \log \mathbb{E}[\exp(\xi W)]$. To construct the new probability measure, we first choose a $\xi_0 > 0$ such that $\psi'(\xi_0) = b$. The new probability measure dF_{ξ_0} is constructed using exponential embedding, as follows

$$dF_{\xi_0} = \exp(\xi_0 W - \psi(\xi_0)) dF,$$

where dF is the original distribution of W . Let \mathbb{E}_{ξ_0} and \mathbb{P}_{ξ_0} denote the expectation and probability under the new measure dF_{ξ_0} , respectively. It can be verified that under the new measure

$$\mathbb{E}_{\xi_0}[W] = \mathbb{E}[W \exp(\xi_0 W - \psi(\xi_0))] = e^{-\psi(\xi_0)} \frac{\partial e^{\psi(\xi)}}{\partial \xi} \Big|_{\xi=\xi_0} = \psi'(\xi_0) = b.$$

Namely, the mean of W is close to the threshold b under the new probability measure.

The threshold crossing probability can be rewritten as

$$\begin{aligned}
\mathbb{P}\left(W = b + \frac{x}{b}\right) &= \mathbb{E}_{\xi_0} \left[\frac{1}{\exp[\xi_0 W - \psi(\xi_0)]} \mathbb{1}\left\{W = b + \frac{x}{b}\right\} \right] \\
&= \exp\left[\psi(\xi_0) - \xi_0 \left(b + \frac{x}{b}\right)\right] \mathbb{P}_{\xi_0}\left(W = b + \frac{x}{b}\right).
\end{aligned} \tag{76}$$

Now we apply the Gaussian approximation to obtain $\mathbb{P}_{\xi_0}\left(W = b + \frac{x}{b}\right)$ and use (76) to get the original probability. By treating W as a normal random variable with mean b and variance $\sigma_{\xi_0}^2$, we have

$$\mathbb{P}_{\xi_0}\left(W = b + \frac{x}{b}\right) = \frac{1}{\sqrt{2\pi}\sigma_{\xi_0}} \exp\left(\frac{-x^2}{2b^2\sigma_{\xi_0}^2}\right) \approx \frac{1}{\sqrt{2\pi}\sigma_{\xi_0}}.$$

Note that in (75), the integrands with smaller x values contribute more to the integration, since the integrand decays exponentially fast with x . Now, when $b \rightarrow \infty$, $\frac{x}{b} \rightarrow 0$ for small x , and hence $\exp\left(\frac{-x^2}{2b^2\sigma_{\xi_0}^2}\right) \rightarrow 1$.

The cumulant generating function of W is calculated as

$$\psi(\xi) = -\xi \frac{\text{tr}(\mathbf{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta))}{\left[2\text{tr}(\mathbf{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta) \mathbf{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta))\right]^{1/2}} - \frac{1}{2} \log \left| \mathbf{I}_{p\tau} - \frac{2\xi \mathbf{\Sigma}_\tau^{1/2} \mathbf{V}_\tau(\theta) \mathbf{\Sigma}_\tau^{1/2}}{\left[2\text{tr}(\mathbf{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta) \mathbf{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta))\right]^{1/2}} \right|.$$

Hence ξ_0 can be obtained by solving the following equation numerically,

$$\frac{1}{\sqrt{d(\tau, \theta)}} \text{tr} \left(\left[\mathbf{I}_{p\tau} - \frac{2\xi_0 \mathbf{B}_\tau(\theta)}{\sqrt{d(\tau, \theta)}} \right]^{-1} \mathbf{B}_\tau(\theta) - \mathbf{A}_\tau(\theta) \right) = b.$$

Eventually, we have

$$\mathbb{P}\left(W\left(i_0, j_0 \frac{\Delta}{\sqrt{N}}\right) = b + \frac{x}{b}\right) \approx g\left(i_0, j_0\right) \exp\left(-\frac{\xi_0}{b}x\right), \quad (77)$$

where $g(\cdot)$ follows the definition in (37).

Step 3: Next we tackle with the conditional probability

$$\mathbb{P}\left(\max_{(i,j) \in J(i_0, j_0)} W\left(i, j \frac{\Delta}{\sqrt{N}}\right) < b \middle| W\left(i_0, j_0 \frac{\Delta}{\sqrt{N}}\right) = b + \frac{x}{b}\right).$$

The first order expansion of the covariance function given by Lemma B.4.2 does not have any cross product term, which implies that if we approximate $W(\tau, \theta)$ as a Gaussian random field, it can be decomposed as a sum of two independent one dimensional random processes.

By Lemma B.4.3, the conditional probability can be written in terms of the decomposed random processes using the techniques in [54] and [24] as follows,

$$\begin{aligned} & \mathbb{P}\left(\max_{(i,j) \in J(i_0, j_0)} W\left(i, j \frac{\Delta}{\sqrt{N}}\right) < b \middle| W\left(i_0, j_0 \frac{\Delta}{\sqrt{N}}\right) = b + \frac{x}{b}\right) \\ &= \mathbb{P}\left(\max_{(i,j) \in J(i_0, j_0)} b \left[W\left(i, j \frac{\Delta}{\sqrt{N}}\right) - W\left(i_0, j_0 \frac{\Delta}{\sqrt{N}}\right) \right] \leq -x \middle| W\left(i_0, j_0 \frac{\Delta}{\sqrt{N}}\right) = b + \frac{x}{b}\right) \\ &\approx \mathbb{P}\left(\max_{i \geq 1} S_i \leq -x\right) \mathbb{P}\left(\max_{i \leq 0} S_i + \max_{j \geq 1} V_j \leq -x\right). \end{aligned} \quad (78)$$

Combine the approximations in (77) and (78), the approximated false alarm rate becomes,

$$\begin{aligned} & \mathbb{P}\left(\max_{(i,j) \in D} W\left(i, j \frac{\Delta}{\sqrt{N}}\right) \geq b\right) \\ & \approx \sum_{(i_0, j_0) \in D} g\left(i_0, j_0 \frac{\Delta}{\sqrt{N}}\right) \frac{\Delta}{\sqrt{N}} \frac{\sqrt{N}}{\Delta b} \int_0^\infty \exp\left(-\frac{\xi_0}{b} x\right) \\ & \quad \times \mathbb{P}\left(\max_{i \geq 1} S_i \leq -x\right) \mathbb{P}\left(\max_{i \leq 0} S_i + \max_{j \geq 1} V_j \leq -x\right) dx. \end{aligned} \quad (79)$$

Finally, by Lemma B.4.4 with $\alpha = \frac{\xi_0}{b}$, $\beta = \sqrt{2}\gamma(\tau, \theta) \frac{b}{\sqrt{N}}$, $\mu_1 = \frac{\mu(\tau, \theta)}{2\tau} b^2$ and $\sigma_1^2 = \frac{\mu(\tau, \theta)}{\tau} b^2$, we have the approximated significance level

$$\frac{1}{2\sqrt{\pi}} \sum_{(i_0, j_0) \in D} g\left(i_0, j_0 \frac{\Delta}{\sqrt{N}}\right) \frac{b^2 \mu(i_0, j_0 \frac{\Delta}{\sqrt{N}})}{N - i_0} \cdot \nu\left(\sqrt{\frac{b^2 \mu(i_0, j_0 \frac{\Delta}{\sqrt{N}})}{N - i_0}}\right) \gamma\left(i_0, j_0 \frac{\Delta}{\sqrt{N}}\right) \frac{\Delta}{\sqrt{N}}. \quad (80)$$

As $\Delta \rightarrow 0$, the Riemann sum (80) converges to the approximation in Theorem 3.3.1.

B.5 Proof of Lemma B.4.1: covariance function of W .

Proof. Let $\mathbf{C}_\tau(\theta) = \mathbf{\Sigma}_\tau^{-1} \mathbf{V}_\tau(\theta) \mathbf{\Sigma}_\tau^{-1}$, and rewrite $\mathbf{C}_m(\theta_2)$ as,

$$\mathbf{C}_m(\theta_2) = \begin{bmatrix} \mathbf{C}_{11}(\theta_2) & \mathbf{C}_{12}(\theta_2) \\ \mathbf{C}_{21}(\theta_2) & \mathbf{C}_n(\theta_2) \end{bmatrix}.$$

Denote $\mathbf{y}_{(T-\tau+1:T)}$ as Y_τ , and let

$$Y_m = \begin{bmatrix} Y_\Delta \\ Y_n \end{bmatrix}.$$

We have,

$$\text{Cov}[W(n, \theta_1), W(m, \theta_2)] = \frac{\mathbb{E}[Y_n^\top \mathbf{C}_n(\theta_1) Y_n Y_m^\top \mathbf{C}_m(\theta_2) Y_m] - \mathbb{E}[Y_n^\top \mathbf{C}_n(\theta_1) Y_n] \mathbb{E}[Y_m^\top \mathbf{C}_m(\theta_2) Y_m]}{2(\text{tr}\{\mathbf{A}_n(\theta_1) \mathbf{A}_n(\theta_1)\} \text{tr}\{\mathbf{A}_m(\theta_2) \mathbf{A}_m(\theta_2)\})^{1/2}}. \quad (81)$$

We calculate the first term in the numerator in the following,

$$\begin{aligned} & \mathbb{E}[Y_n^\top \mathbf{C}_n(\theta_1) Y_n Y_m^\top \mathbf{C}_m(\theta_2) Y_m] \\ & = \mathbb{E}[(Y_n^\top \mathbf{C}_n(\theta_1) Y_n)(Y_\Delta^\top \mathbf{C}_{11}(\theta_2) Y_\Delta + Y_n^\top \mathbf{C}_n(\theta_2) Y_n + Y_n^\top \mathbf{C}_{21}(\theta_2) Y_\Delta + Y_\Delta^\top \mathbf{C}_{12}(\theta_2) Y_n)] \\ & = \mathbb{E}[Y_n^\top \mathbf{C}_n(\theta_1) Y_n Y_n^\top \mathbf{C}_n(\theta_2) Y_n] + \mathbb{E}[Y_n^\top \mathbf{C}_n(\theta_1) Y_n] \mathbb{E}[Y_\Delta^\top \mathbf{C}_{11}(\theta_2) Y_\Delta] \\ & = 2\text{tr}\{\mathbf{A}_n(\theta_1) \mathbf{A}_n(\theta_2)\} + \text{tr}\{\mathbf{A}_n(\theta_1)\} \text{tr}\{\mathbf{A}_n(\theta_2)\} + \mathbb{E}[Y_n^\top \mathbf{C}_n(\theta_1) Y_n] \mathbb{E}[Y_\Delta^\top \mathbf{C}_{11}(\theta_2) Y_\Delta]. \end{aligned} \quad (82)$$

Note that we utilize the fact that under the null hypothesis, Y_Δ and Y_n are independent and $E[Y_\Delta] = 0$. The second term in the numerator is calculated as follows,

$$\begin{aligned}
& E[Y_n^\top \mathbf{C}_n(\theta_1) Y_n] E[Y_m^\top \mathbf{C}_m(\theta_2) Y_m] \\
&= E[Y_n^\top \mathbf{C}_n(\theta_1) Y_n] E[Y_\Delta^\top \mathbf{C}_{11}(\theta_2) Y_\Delta + Y_n^\top \mathbf{C}_n(\theta_2) Y_n + Y_n^\top \mathbf{C}_{21}(\theta_2) Y_\Delta + Y_\Delta^\top \mathbf{C}_{12}(\theta_2) Y_n] \\
&= E[Y_n^\top \mathbf{C}_n(\theta_1) Y_n] E[Y_n^\top \mathbf{C}_n(\theta_2) Y_n] + E[Y_n^\top \mathbf{C}_n(\theta_1) Y_n] E[Y_\Delta^\top \mathbf{C}_{11}(\theta_1) Y_\Delta] \\
&= \text{tr}\{\mathbf{A}_n(\theta_1)\} \text{tr}\{\mathbf{A}_n(\theta_2)\} + E[Y_n^\top \mathbf{C}_n(\theta_1) Y_n] E[Y_\Delta^\top \mathbf{C}_{11}(\theta_1) Y_\Delta].
\end{aligned} \tag{83}$$

By combining (81), (82), (83), we obtain the covariance function in Lemma B.4.1. □

B.6 Proof of Lemma B.4.2: first-order expansion of the covariance function of W .

Proof. We approximate the covariance function by expanding each term in (71) at θ and keeping only the first order terms.

The numerator in (71) is approximated as,

$$\begin{aligned}
\text{tr}(\mathbf{A}_\tau(\theta + \delta) \mathbf{A}_\tau(\theta)) &\approx \text{tr}(\mathbf{A}_\tau(\theta) \mathbf{A}_\tau(\theta)) + \delta \text{tr}(\dot{\mathbf{A}}_\tau(\theta) \mathbf{A}_\tau(\theta)) \\
&= \text{tr}(\mathbf{A}_\tau(\theta) \mathbf{A}_\tau(\theta)) (1 + \delta \gamma(\tau, \theta)).
\end{aligned} \tag{84}$$

We partition the matrix $\mathbf{A}_{\tau+i}(\theta + \delta)$ as follows,

$$\mathbf{A}_{\tau+i}(\theta + \delta) = \begin{bmatrix} \mathbf{A}_{11}(\theta + \delta) & \mathbf{A}_{12}(\theta + \delta) \\ \mathbf{A}_{21}(\theta + \delta) & \mathbf{A}_\tau(\theta + \delta) \end{bmatrix},$$

and rewrite the second term in the denominator in (71) as,

$$\begin{aligned}
& \text{tr}(\mathbf{A}_{\tau+i}(\theta + \delta) \mathbf{A}_{\tau+i}(\theta + \delta)) \\
&= \text{tr}(\mathbf{A}_{11}(\theta + \delta) \mathbf{A}_{11}(\theta + \delta)) + \text{tr}(\mathbf{A}_{12}(\theta + \delta) \mathbf{A}_{21}(\theta + \delta)) \\
&+ \text{tr}(\mathbf{A}_{21}(\theta + \delta) \mathbf{A}_{12}(\theta + \delta)) + \text{tr}(\mathbf{A}_\tau(\theta + \delta) \mathbf{A}_\tau(\theta + \delta)).
\end{aligned}$$

After expanding each term at θ , the denominator in (71) can be approximated as,

$$\left[\text{tr}(\mathbf{A}_\tau(\theta) \mathbf{A}_\tau(\theta)) \text{tr}(\mathbf{A}_{\tau+i}(\theta + \delta) \mathbf{A}_{\tau+i}(\theta + \delta)) \right]^{1/2} \approx \text{tr}(\mathbf{A}_\tau(\theta) \mathbf{A}_\tau(\theta)) \sqrt{1 + 2\delta a} \sqrt{1 + b}, \tag{85}$$

where

$$a = \frac{\text{tr}(\dot{\mathbf{A}}_{11}(\theta) \mathbf{A}_{11}(\theta)) + \text{tr}(\dot{\mathbf{A}}_{12}(\theta) \mathbf{A}_{21}(\theta)) + \text{tr}(\dot{\mathbf{A}}_{21}(\theta) \mathbf{A}_{12}(\theta)) + \text{tr}(\dot{\mathbf{A}}_\tau(\theta) \mathbf{A}_\tau(\theta))}{\text{tr}(\mathbf{A}_{11}(\theta) \mathbf{A}_{11}(\theta)) + \text{tr}(\mathbf{A}_{12}(\theta) \mathbf{A}_{21}(\theta)) + \text{tr}(\mathbf{A}_{21}(\theta) \mathbf{A}_{12}(\theta)) + \text{tr}(\mathbf{A}_\tau(\theta) \mathbf{A}_\tau(\theta))}, \tag{86}$$

and

$$b = \frac{2i}{\tau} \frac{\frac{1}{2i\tau} [\text{tr}(\mathbf{A}_{\tau+i}(\theta)\mathbf{A}_{\tau+i}(\theta)) - \text{tr}(\mathbf{A}_\tau(\theta)\mathbf{A}_\tau(\theta))]}{\frac{1}{\tau^2} \text{tr}(\mathbf{A}_\tau(\theta)\mathbf{A}_\tau(\theta))}. \quad (87)$$

$\dot{\mathbf{A}}(\theta)$ denotes the derivative of a matrix \mathbf{A} with respect to the parameter θ . As i and δ are small compared to τ and θ , the terms $\text{tr}(\dot{\mathbf{A}}_\tau(\theta)\mathbf{A}_\tau(\theta))$ and $\text{tr}(\mathbf{A}_\tau(\theta)\mathbf{A}_\tau(\theta))$ are larger than the subdiagonal elements in (86), and hence, a can be further approximated as,

$$a \approx \frac{\text{tr}(\dot{\mathbf{A}}_\tau(\theta)\mathbf{A}_\tau(\theta))}{\text{tr}(\mathbf{A}_\tau(\theta)\mathbf{A}_\tau(\theta))}. \quad (88)$$

Meanwhile, we approximate the term $\frac{1}{i\tau} [\text{tr}(\mathbf{A}_{\tau+i}(\theta)\mathbf{A}_{\tau+i}(\theta)) - \text{tr}(\mathbf{A}_\tau(\theta)\mathbf{A}_\tau(\theta))]$ in (87) using $\frac{1}{\tau} [\text{tr}(\mathbf{A}_{\tau+1}(\theta)\mathbf{A}_{\tau+1}(\theta)) - \text{tr}(\mathbf{A}_\tau(\theta)\mathbf{A}_\tau(\theta))]$, and then obtain,

$$b \approx \frac{i}{\tau} \mu(\tau, \theta). \quad (89)$$

The argument for the above approximation is as follows. First note that

$$\mathbf{A}_{\tau+i}(\theta) = \Sigma_{\tau+i}^{-1} \mathbf{V}_{\tau+i}(\theta) = (\mathbf{I}_{\tau+i} \otimes \Sigma)^{-1} (\mathbf{R}_{\tau+i}(\theta) \otimes \Lambda) = \mathbf{R}_{\tau+i}(\theta) \otimes (\Sigma^{-1} \Lambda).$$

Then we have,

$$\begin{aligned} \text{tr}(\mathbf{A}_{\tau+i}(\theta)\mathbf{A}_{\tau+i}(\theta)) &= \text{tr}((\mathbf{R}_{\tau+i}(\theta) \otimes (\Sigma^{-1} \Lambda))(\mathbf{R}_{\tau+i}(\theta) \otimes (\Sigma^{-1} \Lambda))) \\ &= \text{tr}((\mathbf{R}_{\tau+i}(\theta)\mathbf{R}_{\tau+i}(\theta)) \otimes (\Sigma^{-1} \Lambda \Sigma^{-1} \Lambda)) \\ &= \text{tr}(\mathbf{R}_{\tau+i}(\theta)\mathbf{R}_{\tau+i}(\theta)) \text{tr}(\Sigma^{-1} \Lambda \Sigma^{-1} \Lambda) \\ &= \text{tr}(\Sigma^{-1} \Lambda \Sigma^{-1} \Lambda) \sum_j \sum_k [\mathbf{R}_{\tau+i}(\theta)]_{jk}^2 \\ &= \text{tr}(\Sigma^{-1} \Lambda \Sigma^{-1} \Lambda) \left(i \sum_j \sum_k [\mathbf{R}_{\tau+1}(\theta)]_{jk}^2 + \sum_{|j-k|>\tau} [\mathbf{R}_{\tau+1}(\theta)]_{jk}^2 \right) \\ &\approx \text{tr}(\Sigma^{-1} \Lambda \Sigma^{-1} \Lambda) \left(i \sum_j \sum_k [\mathbf{R}_{\tau+1}(\theta)]_{jk}^2 \right). \end{aligned}$$

The last approximation is due to the fact that (j, k) th element of $\mathbf{R}_{\tau+1}(\theta)$ such that $|j-k| > \tau$ is small. Combining (84), (85), (88) (89) and the Taylor expansion $\frac{1}{\sqrt{1+x}} \approx 1 - \frac{1}{2}x + o(x)$, we obtain the approximation in (72).

□

REFERENCES

- [1] BARTRAM, J. and BALLANCE, R., *Water quality monitoring: a practical guide to the design and implementation of freshwater quality studies and monitoring programmes*. CRC Press, 1996.
- [2] BERNARDO, J., BAYARRI, M., BERGER, J., DAWID, A., HECKERMAN, D., SMITH, A., and WEST, M., “Optimization under unknown constraints,” *Bayesian Statistics*, vol. 9, no. 9, p. 229, 2011.
- [3] BODNAR, O. and SCHMID, W., “Multivariate control charts based on a projection approach,” *Allgemeines Statistisches Archiv*, vol. 89, pp. 75–93, 3005.
- [4] BONILLA, E. V., CHAI, K. M., and WILLIAMS, C., “Multi-task gaussian process prediction,” in *Advances in neural information processing systems*, pp. 153–160, 2008.
- [5] BROCKWELL, P. J. and DAVIS, R. A., *Time Series: Theory and Methods*. Springer-Verlag New York, 1987.
- [6] BRYNJARSDÓTTIR, J. and BERLINER, L. M., “Dimension-reduced modeling of spatio-temporal processes,” *Journal of the American Statistical Association*, vol. 109, no. 508, pp. 1647–1659, 2014.
- [7] CHEN, J., KIM, S.-H., and XIE, Y., “S³T: A score statistic for spatio-temporal change-point detection,” *arXiv preprint*, 2018.
- [8] CLEMENT, L. and THAS, O., “Estimating and modeling spatio-temporal correlation structures for river monitoring networks,” *Journal of Agricultural, Biological, and Environmental Statistics*, vol. 12, pp. 161–176, Jun 2007.
- [9] CLEMENT, L., THAS, O., VANROLLEGHEM, P., and OTTOY, J., “Spatio-temporal statistical models for river monitoring networks,” *Water Science and Technology*, vol. 53, no. 1, pp. 9–15, 2006.
- [10] CORLESS, R. M., GONNET, G. H., HARE, D. E., JEFFREY, D. J., and KNUTH, D. E., “On the lambertw function,” *Advances in Computational mathematics*, vol. 5, no. 1, pp. 329–359, 1996.
- [11] CROSIER, R. B., “Multivariate generalizations of cumulative sum quality-control schemes,” *Technometrics*, vol. 30, no. 3, pp. 291–303, 1988.
- [12] GAETAN, C. and GUYON, X., *Spatial Statistics and Modeling*. Springer, New York, 2010.
- [13] GAETAN, C. and GUYON, X., *Spatial statistics and modeling*, vol. 81. Springer, 2010.
- [14] GARDNER, J. R., KUSNER, M. J., XU, Z. E., WEINBERGER, K. Q., and CUNNINGHAM, J. P., “Bayesian optimization with inequality constraints,” in *ICML*, pp. 937–945, 2014.

- [15] GARNETT, R., OSBORNE, M. A., and ROBERTS, S. J., “Bayesian optimization for sensor set selection,” in *Proceedings of the 9th ACM/IEEE international conference on information processing in sensor networks*, pp. 209–219, ACM, 2010.
- [16] GENTON, M. G., “Separable approximations of space-time covariance matrices,” *Environmetrics*, vol. 18, no. 7, pp. 681–695, 2007.
- [17] GUERRIERO, M., WILLETT, P., and GLAZ, J., “Distributed target detection in sensor networks using scan statistics,” *Signal Processing, IEEE Transactions on*, vol. 57, no. 7, pp. 2629–2639, 2009.
- [18] HE, Q. and ZHOU, S., “Discriminant locality preserving projection chart for statistical monitoring of manufacturing processes,” *International Journal of Production Research*, vol. 52, no. 18, pp. 5286–5300, 2014.
- [19] HEALY, J. D., “A note on multivariate cusum procedures,” *Technometrics*, vol. 29, no. 4, pp. 409–412, 1987.
- [20] HERNÁNDEZ-LOBATO, J. M., GELBART, M. A., ADAMS, R. P., HOFFMAN, M. W., and GHAHRAMANI, Z., “A general framework for constrained bayesian optimization using information-based search,” *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 5549–5601, 2016.
- [21] HOTELLING, H., “Multivariate quality control,” *Techniques of statistical analysis*, 1947.
- [22] JIANG, W., HAN, S. W., TSUI, K.-L., and WOODALL, W. H., “Spatiotemporal surveillance methods in the presence of spatial correlation,” *Statistics in Medicine*, vol. 30, no. 5, pp. 569–583, 2011.
- [23] JONES, D. R., SCHONLAU, M., and WELCH, W. J., “Efficient global optimization of expensive black-box functions,” *Journal of Global optimization*, vol. 13, no. 4, pp. 455–492, 1998.
- [24] KIM, H.-J. and SIEGMUND, D., “The likelihood ratio test for a change-point in simple linear regression,” *Biometrika*, vol. 76, pp. 409–423, 1989.
- [25] KIM, S.-H., , ALEXOPOULOS, C., TSUI, K.-L., and WILSON, J. R., “A distribution-free tabular CUSUM chart for autocorrelated data,” *IIE Transactions*, vol. 39, no. 3, pp. 317–330, 2007.
- [26] KIM, S.-H., ARAL, M. M., EUN, Y., PARK, J. J., and PARK, C., “Impact of sensor measurement error on sensor positioning in water quality monitoring networks,” *Stochastic Environmental Research and Risk Assessment*, vol. 31, no. 3, pp. 743–756, 2017.
- [27] LAM, R. and WILLCOX, K., “Lookahead bayesian optimization with inequality constraints,” in *Advances in Neural Information Processing Systems*, pp. 1890–1900, 2017.
- [28] LEE, J., HUR, Y., KIM, S.-H., and WILSON, J. R., “Monitoring nonlinear profiles using a wavelet-based distribution-free cusum chart,” *International Journal of Production Research*, vol. 50, no. 22, pp. 6574–6594, 2012.

- [29] LEE, M. L., GOLDSMAN, D., and KIM, S.-H., “Robust distribution-free multivariate cusum charts for spatiotemporal biosurveillance in the presence of spatial correlation,” *IIE Transactions on Healthcare Systems Engineering*, vol. 5, no. 2, pp. 74–88, 2015.
- [30] LEE, M. L., GOLDSMAN, D., KIM, S.-H., and TSUI, K.-L., “Spatiotemporal biosurveillance with spatial clusters: control limit approximation and impact of spatial correlation,” *IIE Transactions*, vol. 46, no. 8, pp. 813–827, 2014.
- [31] LETHAM, B., KARRER, B., OTTONI, G., BAKSHY, E., and OTHERS, “Constrained bayesian optimization with noisy experiments,” *Bayesian Analysis*, 2018.
- [32] LIU, K., ZHANG, R., and MEI, Y., “Scalable sum-shrinkage schemes for distributed monitoring large-scale data streams,” *arXiv:1603.08652*, 2016.
- [33] LIU, K., MEI, Y., and SHI, J., “An adaptive sampling strategy for online high-dimensional process monitoring,” *Technometrics*, vol. 57, no. 3, pp. 305–319, 2015.
- [34] LOWRY, C. A., WOODALL, W. H., CHAMP, C. W., and RIGDON, S. E., “A multivariate exponentially weighted moving average control chart,” *Technometrics*, vol. 34, no. 1, pp. 46–53, 1992.
- [35] MISHIN, D., BRANTNER-MAGEE, K., CZAKO, F., and SZALAY, A. S., “Real time change point detection by incremental pca in large scale sensor data,” in *High Performance Extreme Computing Conference (HPEC), 2014 IEEE*, pp. 1–6, IEEE, 2014.
- [36] MOCKUS, J., TIESIS, V., and ZILINSKAS, A., “The application of bayesian methods for seeking the extremum,” vol. 2, pp. 117–129, 1978.
- [37] NASA, “SDO instruments,” Retrieved 7-30-2012.
- [38] PAGE, E. S., “Continuous inspection schemes,” *Biometrika*, vol. 41, no. 1/2, pp. 100–115, 1954.
- [39] PARK, C. and KIM, S.-H., “Penalty function with memory for discrete optimization via simulation with stochastic constraints,” *Operations Research*, vol. 63, no. 5, pp. 1195–1212, 2015.
- [40] PARK, C., KIM, S.-H., TELCI, I. T., and ARAL, M. M., “Designing optimal water quality monitoring network for river systems and application to a hypothetical river,” in *Proceedings of the Winter Simulation Conference*, pp. 3506–3513, Winter Simulation Conference, 2010.
- [41] PARK, C., TELCI, I. T., KIM, S.-H., and ARAL, M. M., “Designing an optimal water quality monitoring network for river systems using constrained discrete optimization via simulation,” *Engineering Optimization*, vol. 46, no. 1, pp. 107–129, 2014.
- [42] PARK, Y., BAEK, S. H., KIM, S.-H., and TSUI, K.-L., “Statistical process control-based intrusion detection and monitoring,” *Quality and Reliability Engineering International*, vol. 30, no. 2, pp. 257–273, 2014.
- [43] RAO, C. R., “Large sample tests of statistical hypotheses concerning several parameters with applications to problems of estimation,” *Mathematical Proceedings of the Cambridge Philosophical Society*, vol. 44, no. 1, pp. 50–57, 1948.

- [44] RAO, C. R. and POTI, S. J., “On locally most powerful tests when alternatives are one sided,” *Sankhyā: The Indian Journal of Statistics*, vol. 7, pp. 439–439, 1946.
- [45] RASMUSSEN, C. E. and WILLIAMS, C. K. I., *Gaussian processes for machine learning*. MIT Press, Cambridge, MA, 2006.
- [46] RIPLEY, B. D., *Spatial statistics*, vol. 575. John Wiley & Sons, 2005.
- [47] ROGERSON, P. A., “Surveillance systems for monitoring the development of spatial patterns,” *Statistics in Medicine*, vol. 16, no. 18, pp. 2081–2093, 1997.
- [48] ROGERSON, P. A. and YAMADA, I., “Approaches to syndromic surveillance when data consist of small regional counts,” *Morbidity and Mortality Weekly Report*, pp. 79–85, 2004.
- [49] ROSSMAN, L. A., *Storm water management model user’s manual, version 5.0*. National Risk Management Research Laboratory, Office of Research and Development, US Environmental Protection Agency, 2010.
- [50] RUBNER, Y., TOMASI, C., and GUIBAS, L. J., “The earth mover’s distance as a metric for image retrieval,” *International journal of computer vision*, vol. 40, no. 2, pp. 99–121, 2000.
- [51] RUNGER, G. C., “Projections and the U-squared multivariate control chart,” *J. of Quality Technology*, vol. 28, no. 3, pp. 313–319, 1996.
- [52] SHI, L. and OTHERS, “Nested partitions method for stochastic optimization,” *Methodology and Computing in Applied probability*, vol. 2, no. 3, pp. 271–291, 2000.
- [53] SIEGMUND, D. and YAKIR, B., *The statistics of gene mapping*. Springer, 2007.
- [54] SIEGMUND, D., “Approximate tail probabilities for the maxima of some random fields,” *The Annals of Probability*, vol. 16, pp. 487–501, 1988.
- [55] SIEGMUND, D., *Sequential analysis: tests and confidence intervals*. Springer Science & Business Media, 2013.
- [56] SIEGMUND, D. and VENKATRAMAN, E., “Using the generalized likelihood ratio statistic for sequential detection of a change-point,” *The Annals of Statistics*, vol. 23, pp. 255–271, 1995.
- [57] SIEGMUND, D. and YAKIR, B., “Detecting the emergence of a signal in a noisy image,” *Statistics and Its Inference*, vol. 1, pp. 3–12, 2008.
- [58] SKUBALLSKA-RAFAJLOWICZ, E., “Random projections and Hotelling’s t^2 statistics for change detection in high-dimensional data analysis,” *Int. J. Appl. Math. Comput. Sci.*, vol. 23, no. 2, pp. 447–461, 2013.
- [59] SNOEK, J., LAROCHELLE, H., and ADAMS, R. P., “Practical bayesian optimization of machine learning algorithms,” in *Advances in neural information processing systems*, pp. 2951–2959, 2012.
- [60] SONESSON, C., “A cusum framework for detection of space–time disease clusters using scan statistics,” *Statistics in Medicine*, vol. 26, no. 26, pp. 4770–4789, 2007.

- [61] SPIEGELHALTER, D., SHERLAW-JOHNSON, C., BARDSLEY, M., BLUNT, I., WOOD, C., and GRIGG, O., “Statistical methods for healthcare regulation: rating, screening and surveillance,” *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, vol. 175, no. 1, pp. 1–47, 2012.
- [62] SRINIVAS, N., KRAUSE, A., KAKADE, S. M., and SEEGER, M., “Gaussian process optimization in the bandit setting: No regret and experimental design,” *arXiv preprint arXiv:0912.3995*, 2009.
- [63] STROBL, R. O. and ROBILLARD, P. D., “Network design for water quality monitoring of surface freshwaters: A review,” *Journal of environmental management*, vol. 87, no. 4, pp. 639–648, 2008.
- [64] TANGO, T., “A class of tests for detecting ‘general’ and ‘focused’ clustering of rare diseases,” *Statistics in Medicine*, vol. 14, no. 21-22, pp. 2323–2334, 1995.
- [65] TELCI, I. T. and ARAL, M. M., “Contaminant source location identification in river networks using water quality monitoring systems for exposure analysis,” *Water Quality, Exposure and Health*, vol. 2, no. 3-4, pp. 205–218, 2011.
- [66] TELCI, I. T., NAM, K., GUAN, J., and ARAL, M. M., “Real time optimal monitoring network design in river networks,” in *World environmental & water resources congress*, 2008.
- [67] TELCI, I. T., NAM, K., GUAN, J., and ARAL, M. M., “Optimal water quality monitoring network design for river systems,” *Journal of environmental management*, vol. 90, no. 10, pp. 2987–2998, 2009.
- [68] TSUI, K.-L., CHIU, W., GIERLICH, P., GOLDSMAN, D., LIU, X., and MASCHKE, T., “A review of healthcare, public health, and syndromic surveillance,” *Quality Engineering*, vol. 20, no. 4, pp. 435–450, 2008.
- [69] VER HOEF, J. M. and PETERSON, E. E., “A moving average approach for spatial statistical models of stream networks,” *Journal of the American Statistical Association*, vol. 105, no. 489, pp. 6–18, 2010.
- [70] WANG, H., KIM, S.-H., HUO, X., HUR, Y., and WILSON, J. R., “Monitoring non-linear profiles adaptively with a wavelet-based distribution-free cusum chart,” *International Journal of Production Research*, vol. 53, no. 15, pp. 4648–4667, 2015.
- [71] WANG, K., JIANG, W., and LI, B., “A spatial variable selection method for monitoring product surface,” *International Journal of Production Research*, vol. 54, no. 14, pp. 4161–4181, 2015.
- [72] XIE, Y., HUANG, J., and WILLETT, R., “Change-point detection for high-dimensional time series with missing data,” *IEEE Journal of Selected Topics in Signal Processing (J-STSP)*, vol. 7, pp. 12–27, Feb. 2013.
- [73] XIE, Y. and SIEGMUND, D., “Spectrum opportunity detection with weak and correlated signals,” in *2012 Conference Record of the Forty Sixth Asilomar Conference on Signals, Systems and Computers (ASILOMAR)*, pp. 128–132, IEEE, 2012.

- [74] XIE, Y. and SIEGMUND, D., “Sequential multi-sensor change-point detection,” *The Annals of Statistics*, vol. 41, no. 2, pp. 670–692, 2013.
- [75] XIE, Y., WANG, M., and THOMPSON, A., “Sketching for sequential change-point detection,” in *2015 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, pp. 78–82, IEEE, 2015.
- [76] YAKIR, B., *Extremes in random fields: A theory and its applications*. John Wiley & Sons, 2013.
- [77] YAMAMOTO, M. and HWANG, H., “Dimension-reduced clustering of functional data via subspace separation,” *Journal of Classification*, vol. 34, no. 2, pp. 294–326, 2017.
- [78] ZHANG, F., *The Schur complement and its applications*, vol. 4. Springer Science & Business Media, 2006.